

HALLUSIONBENCH: An Advanced Diagnostic Suite for Entangled Language Hallucination and Visual Illusion in Large Vision-Language Models

Tianrui Guan* Fuxiao Liu* Xiyang Wu Ruiqi Xian Zongxia Li Xiaoyu Liu Xijun Wang
Lichang Chen Furong Huang Yaser Yacoob Dinesh Manocha Tianyi Zhou

University of Maryland, College Park

{rayguan, fl3es, wuxiyang, rxian, zli12321, xliu1231, xijun
bobchen, furongh, yaser, dmanocha, tianyi}@umd.edu

Abstract

We introduce “HALLUSIONBENCH¹,” a comprehensive benchmark designed for the evaluation of image-context reasoning. This benchmark presents significant challenges to advanced large visual-language models (LVLMs), such as GPT-4V(ision), Gemini Pro Vision, Claude 3, and LLaVA-1.5, by emphasizing nuanced understanding and interpretation of visual data. The benchmark comprises 346 images paired with 1129 questions, all meticulously crafted by human experts. We introduce a novel structure for these visual questions designed to establish control groups. This structure enables us to conduct a quantitative analysis of the models’ response tendencies, logical consistency, and various failure modes. In our evaluation on HALLUSIONBENCH, we benchmarked 15 different models, highlighting a 31.42% question-pair accuracy achieved by the state-of-the-art GPT-4V. Notably, all other evaluated models achieve accuracy below 16%. Moreover, our analysis not only highlights the observed failure modes, including language hallucination and visual illusion but also deepens an understanding of these pitfalls. Our comprehensive case studies within HALLUSIONBENCH shed light on the challenges of hallucination and illusion in LVLMs. Based on these insights, we suggest potential pathways for their future improvement. The benchmark and codebase can be accessed at <https://github.com/tianyi-lab/HallusionBench>.

1. Introduction

In recent years, Large Language Models (LLMs) [8, 9, 25, 40, 45, 46, 61] have revolutionized the field of machine learning with the ability of language understanding and content generation, offering unprecedented ca-

pabilities and potentials across a multitude of applications. The integration of LLMs with computer vision systems has given rise to Large Vision-Language Models (LVLMs) [5, 7, 21, 26, 27, 32, 40, 41, 49, 50, 55, 63]. These models have demonstrated profound capabilities in various applications and significantly enhance the performance in image reasoning tasks [4, 17, 19, 29, 30, 35, 37, 42, 47]. However, the hallucination issue of LLMs [58] is regarded as a challenging and unsolved problem, which leads to many issues when we integrate LLMs with vision techniques.

While LVLMs like GPT-4V(ision) [48] and LLaVA-1.5 [31] excel in various applications, they are hindered by a pronounced language bias. This bias stems from instances where knowledge priors conflict with the visual context [23, 28, 57]. Similarly, models such as LLaVA-1.5 [31] and mPLUG-Owl [50] are prone to giving affirmative answers regardless of the actual content of questions [23]. The distinct failure modes of different VLMs highlight the need for specific improvements. Recognizing and understanding these limitations and failure types is imperative for advancing these models and striking a delicate balance between knowledge priors and contextual understanding.

When exploring those LVLMs, we observe that their strong language bias often overshadows visual information, leading to an overreliance on language priors rather than the visual context. To study this phenomenon, we use the term “**Language Hallucination**,” which refers to conclusions drawn without visual input. On the other hand, the vision components within the limited ability in LVLMs can give rise to “**Visual Illusion**”, where visual inputs can be misinterpreted, leading to overconfident yet erroneous assertions by the model.

Main Contributions: Recognizing the need to comprehend why an LVLM fails and address these issues, we present HALLUSIONBENCH, a carefully crafted benchmark designed to explore the complexities of image-context rea-

*Equal contribution.

¹“Hallusion” is a portmanteau of “hallucination” and “illusion.”

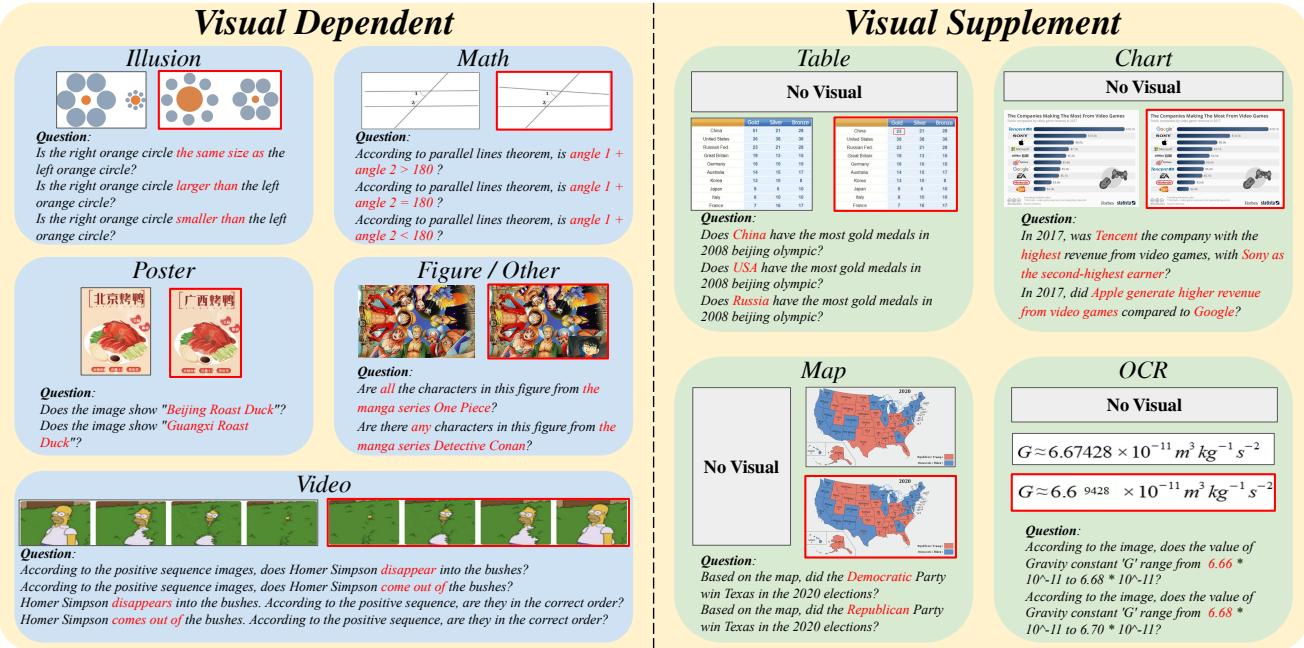


Figure 1. Data samples of HALLUSIONBENCH, which contains diverse topics, visual modalities. Human-edited images are in RED, resulting in different correct answers to the questions.

soning in depth and expose various problems with respect to current LVLMs, as shown in Fig. 1. Our design of the visual-question (VQ) pairs, unique in format, facilitates a quantitative analysis of the models’ failures, enabling a more thorough evaluation. This investigation sheds light on existing limitations and lays the groundwork for future improvements, aiming to make the next generation of LVLMs more robust, balanced, and precise. The novelties of our work include:

1. We introduce HALLUSIONBENCH, the first advanced diagnostic suite tailored to systematically dissect and analyze the diverse failure modes of LVLMs. HALLUSIONBENCH consists of approximately 1129 handcrafted visual question-answer (VQA) pairs, featuring 165 original images and 181 images expertly modified by human professionals. Moving beyond the traditional metrics of correctness and accuracy, our VQA pairs are thoughtfully formulated with an innovative structure. This approach enables us to quantitatively analyze specific dimensions and aspects where current models falter.
2. We evaluate 15 most recent methods on HALLUSIONBENCH. Our benchmark presents formidable challenges to existing methods. Notably, the SoTA GPT-4V achieves merely a 31.42% Question Pair Accuracy, while the performance of all other methods falls below 16%.
3. We explore HALLUSIONBENCH and provide an in-depth analysis of examples on which the SoTA LVLMs, such as GPT-4V and LLaVA-1.5 fail. We also provide insights on different issues that existing LVLMs are facing based on

the quantitative analysis enabled by HALLUSIONBENCH. In our exploration of HALLUSIONBENCH, we conduct a detailed analysis of instances where SoTA LVLMs, including GPT-4V and LLaVA-1.5, fall short. Additionally, our investigation leverages the quantitative capabilities of HALLUSIONBENCH to shed light on various issues currently challenging existing LVLMs.

2. Related Work

2.1. Large Multi-Modal Models

Large Language Models have been a major advancement, leading to new ways to understand not just text but other things like images, all in one large system. For example, Flamingo [3] has many capabilities, combining a vision part that doesn’t change with a big language model that has a special feature for understanding both images and words together. Another model, PaLM-E [12], mixes visual information directly into the already powerful PaLM model, which has 520 billion parameters, making it effective in real-world uses. Most recently, researchers have been creating high-quality, diverse multi-modal datasets from GPT4 and GPT-4V [48] to fine-tune open-source LVLMs, including LLaVA [32], MiniGPT4 [63], Mplug-Owl [50], LRV-Instruction [28], LLaVAR [60] and other works [11, 24, 36, 52].

2.2. Hallucination in LVLMs

Hallucination typically refers to situations where the generated responses contain information that is not present in

the visual content. Prior research primarily examines two areas: detecting and evaluating hallucinations [23, 58, 59], and methods to reduce them [28, 43, 53]. Early methods include training classifiers to identify hallucinations or comparing output with accurate answers to detect inaccuracies. To mitigate hallucinations, efforts have been made to improve data gathering and training procedures. For example, LRV-Instruction [28] creates balanced positive and negative instructions to finetune LVLMs. VIGC [43] uses an iterative process to generate concise answers and combine them, aiming for detailed yet accurate responses. Similarly, Woodpecker [53] introduces a training-free method to pick out and correct hallucinations from the generated text.

2.3. Benchmarks for Large VL Models

Traditional Visual Language (VL) benchmarks are designed to assess distinct skills, including visual recognition [16], image description [2, 27], and so on. However, with the advent of advanced LVLMs, traditional evaluation metrics often fall short of providing a detailed ability assessment. This problem is further exacerbated by their inability to match the given answer accurately, leading to significant robustness issues. To address these challenges, research communities have introduced a series of benchmarks, including MME [14], MMBench [33], MM-Vet [54], SEED-Bench [20], GAVIE [28], and LAMM-Bench [13]. These benchmarks systematically structure and evaluate complex multi-modal tasks. Different from POPE [23] and GAVIE [28] evaluating the object hallucinations of LVLMs, HALLUSIONBENCH is the first human-annotated analytical benchmark focusing on diagnosing both the visual illusion and knowledge hallucination of LVLMs.

3. HALLUSIONBENCH Construction

We present HALLUSIONBENCH, the first benchmark designed to examine visual illusion and knowledge hallucination of LVLMs and analyze the potential failure modes based on each hand-crafted example pair. HALLUSIONBENCH consists of 455 visual-question control pairs, including 346 different figures and a total of 1129 questions on diverse topics (including ***food, math, geometry, statistics, geography, sports, cartoon, famous illusions, movie, meme, etc.***) and formats (including ***logo, poster, figure, charts, table, map, consecutive images, etc.***). In the following sections, we first provide the guidelines for dataset construction based on different visual question types. Second, we will describe the data and annotation structure of HALLUSIONBENCH. Finally, we will describe the statistics of our dataset.

3.1. Visual Question Taxonomy

Our aim is to develop a multimodal image-context reasoning benchmark to investigate the potent language bias inherent in LVLMs, which can sometimes overshadow the visual

context. We define the two categories of visual questions: ***Visual Dependent*** and ***Visual Supplement***.

3.1.1 Visual Dependent Questions

The ***Visual Dependent*** questions are defined as ***questions that do not have an affirmative answer without the visual context***. Such questions ask about the image itself or something within the image. For example, there is no clear answer to "*Is the right orange circle the same size as the left orange circle?*" without an image to provide more context.

Guideline: Under this setting, our benchmark is designed to evaluate visual commonsense knowledge and visual reasoning skills. Our exploration and dataset construction are guided by the following questions:

1. *How good are the visual understanding and reasoning skills of the model?*
2. *How does the parametric memory of the model affect its response to a question?*
3. *Is the model able to capture the temporal relation of multiple images?*

3.1.2 Visual Supplement Questions

The ***Visual Supplement*** questions are ***questions that can be answered without the visual input; the visual component merely provides supplemental information or corrections***. For example, some LVLMs can answer "*Is New Mexico state larger than Texas state?*" using the prior knowledge in their parametric memory without a map of the US.

Guideline: Under this setting, our benchmark is designed to evaluate visual reasoning ability and the balance between parametric memory and image context. Our exploration and dataset construction under this category is guided by the following questions:

1. *When the model lacks the prior knowledge or answer in the parametric memory of its language module, does the model (still) hallucinate about the images?*
2. *When the model's language module has sufficient prior knowledge in its parametric memory or directly knows the answer, does it still enhance its response by gathering extra information from the visual supplement (especially when the prior knowledge conflicts with the visual input or the parametric memory is outdated)?*
3. *How well can the model interpret a visual input with dense information (i.e., a graph, chart, map, etc.) for question answering? What types of image manipulation might impede or distort visual information extraction?*

3.2. Visual, Question, and Annotation Structures

Notations: Let $(I, q) \in \mathcal{V} \subseteq \mathbb{I} \times \mathbb{Q}$ be the tuple of the image $I \in \mathbb{I}$ and question $q \in \mathbb{Q}$, where \mathcal{V} is the set of valid VQ pairs. Let N be the number of original images obtained from the Internet, and $\mathbb{I}_o = \{I_{(i,0)}\}_{0 < i \leq N}$ be the set of those

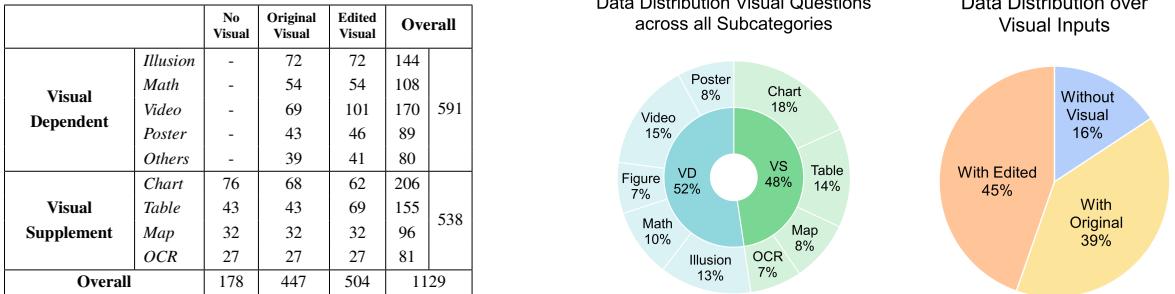


Figure 2. **Statistics of HALLUSIONBENCH:** We show the number of questions in the table (*left*), and the distribution of visual questions across each subcategory of Visual Dependent (VD) and Visual Supplement (VS) (*middle*) and visual input types categorized by no visual, original, and edited images (*right*). HALLUSIONBENCH covers a diverse visual format and nearly half of the images are manually edited.

Benchmarks	Visaul Format	# Total QA	# H-Edited QA	# Total Img.	# H-Edited Img.	Control Pair?	Purpose
Lynx-Bench [56]	Image, Video	450	450	450	0	✗	Image&Video QA Evaluation
SciGraphQA [22]	Image	295K	0	657K	0	✗	Scientific Chart QA Evaluation
MathVista [34]	Image	6141	0	5487	0	✗	Math Reasoning Evaluation
MME [14]	Image	1457	1457	1187	0	✗	Comprehensive Evaluation
POPE [23]	Image	3000	0	500	0	✗	Object Hallucination
M-HalDetect [18]	Image	4000	0	4000	0	✗	Object Hallucination
GAVIE [28]	Image	1000	0	1000	0	✗	Object Hallucination
Bingo [10]	Image	370	370	308	N/A	✓	Hallucination, Bias
HALLUSIONBENCH	Image, Video Image Pairs	1129	1129	346	181	✓	Visual Illusion, Language Hallucination, Quantitative Analysis and Diagnosis

Table 1. **Comparison of HALLUSIONBENCH with most recent VL benchmarks:** HALLUSIONBENCH is the **first** and the **only** benchmark that focuses on control-group analysis by carefully editing each image in the database manually. “# H-Edited QA” means Human-edited question-answer pairs. “# H-Edited Img” means Human-edited images. N/A denotes that the information is not provided.

original images. We define $\mathbb{I}'_i = \{I_{(i,j)}\}_{0 < j \leq N_i}$ be the set of images modified from $I_{(i,0)}$, and I_0 be an empty image. The entire images set $\mathbb{I} = \{I_0\} \cup \mathbb{I}_o \cup (\bigcup_{0 < i \leq N} \mathbb{I}'_i)$.

Let $\mathbb{Q}_i = \{q_{(i,k)}\}_{0 < k \leq M_i}$ be the set of questions that can be applied to any image in \mathbb{I}_i , which is defined differently for Visual Dependent (VD) and Visual Supplement (VS):

$$\mathbb{I}_i = \begin{cases} \{I_{(i,0)}\} \cup \mathbb{I}'_i & \text{for } VD \\ \{I_0, I_{(i,0)}\} \cup \mathbb{I}'_i & \text{for } VS \end{cases} \quad (1)$$

To facilitate evaluation, all questions are formulated as Yes/No questions (Fig. 1). We annotate each visual-question with a binary answer $y(I, q) \in \{\text{"yes"}, \text{"no"}\}$.

3.3. Dataset Statistics

Following the annotation structure and guidelines above, we ask human experts to collect 346 images with diverse topics and types manually. As shown Fig. 2, *Visual Dependent* has 591 questions, including *videos*, *illusion*, *math*, *posters*, *logos*, *cartoons*, and *others*; *Visual Supplement* has 538 questions, including *charts*, *tables*, *maps*, and *OCR*. Furthermore, Fig. 2 (*right*) describes the distribution of the questions without visual input (16%), with original online images (39%), and with visual input edited by human experts (45%). Our image manipulation strategies contain *image flipping*, *order reversing*, *masking*, *optical character editing*, *object editing*, and *color editing*. Additionally, each image has 3.26 questions on average. Fig. 2 (*left*) provides more details on the number of questions in each topic and visual input category.

3.4. Uniqueness of HALLUSIONBENCH

The main comparison between HALLUSIONBENCH and existing benchmarks is presented in Tab. 1. As it shows, there is a notable gap between existing benchmarks [10, 18, 23, 28] and HALLUSIONBENCH in hallucination evaluation, as existing benchmarks primarily focus on object hallucinations, limited topics, and visual input types. Our dataset, HALLUSIONBENCH, is therefore motivated to bridge this gap by providing more topics, more image types, and more visual input modalities, including both images and videos. Additionally, our human experts carefully select each image and write question-answer pairs. We are also the first work to include human-edited images to assess the robustness of current LVLMs. Additionally, unlike existing benchmarks, HALLUSIONBENCH focuses on evaluating both language hallucinations and visual illusions, moving beyond the narrow scope of object hallucinations [18, 23, 28].

4. HALLUSIONBENCH Evaluation Suite

4.1. Text-Only GPT4-Assisted Evaluation

Notations: Let $\mathcal{M}(I, q) \in \{\text{"yes"}, \text{"no"}, \text{"uncertain"}\}$ be the parsed output answer by a VLM \mathcal{M} for an image-question pair (I, q) . GPT-4 $\text{GPT}(\mathcal{M}(I, q), y(I, q))$ then judges the answer $\mathcal{M}(I, q)$ based on the ground truth $y(I, q) \in \{\text{"yes"}, \text{"no"}\}$ and outputs *Incorrect* (0), *Correct* (1), or *Uncertain* (2) if the predicted response is ambiguous.

The prompt for the GPT-4 judge is designed as:

Imagine you are an intelligent teacher. Thoroughly read the question, reference answer, and the prediction answer to ensure a clear understanding of the information provided. Assess the correctness of the predictions. If the prediction answer does not conflict with the reference answer, please generate “correct”. If the prediction answer conflicts with the reference answer, please generate “incorrect”. If the prediction answer is unclear about the answer, please generate “unclear”.

For each sample, we fill the template with its question, ground truth, and LVLM output. By taking the filled prompt into GPT-4, GPT-4 will generate "correct", "incorrect" or "unclear" for the sample. It is found that outputs of GPT-4 still exist variance, although the temperature is set as 0. Therefore, we utilize GPT-4 to evaluate the outputs of LLMs 3 times and report average scores.

Comparison with Human Evaluation: To demonstrate that our GPT4-Assisted evaluation is effective, we obtain the responses from GPT-4V [48] and LLaVA-1.5 [31], and manually evaluate the correctness of their responses. We label the responses with *Incorrect* (0), *Correct* (1), and *Uncertain* (2) if the answer is ambiguous. As shown in the first two rows of Tab. 2 and Tab. 3, the negligible difference proves that the GPT4-assisted method aligns well with human judgment.

4.2. Correctness Evaluation Metrics

Since the focus of our benchmark is on hallucination and illusion, not the span of knowledge, we consider an *uncertain* answer acceptable when there is no visual input under the *Visual Supplement* category. For the final accuracy score, we convert the correctness into a binary value $b_{\mathcal{M}} \in \{0, 1\}$:

$$b_{\mathcal{M}}(I, q) = \begin{cases} GPT(\mathcal{M}(I, q), y(I, q)) & \text{if } GPT(\mathcal{M}, y) \leq 1 \\ 1 & \text{else if } I = I_0 \\ 0 & \text{otherwise} \end{cases}, \quad (2)$$

Let $(I, q) \in \mathcal{V} \subseteq \mathbb{I} \times \mathbb{Q}$ be the tuple of the image $I \in \mathbb{I}$ and question $q \in \mathbb{Q}$, where \mathcal{V} is the set of valid visual-question pairs. Let $\mathbb{1}(\cdot)$ be the indicator function.

All accuracy:

$$aAcc = \frac{\sum_{(I,q) \in \mathcal{V}} b_{\mathcal{M}}(I, q)}{|\mathcal{V}|} \quad (3)$$

Figure Accuracy:

$$fAcc = \frac{\sum_{i,j} \mathbb{1}(\bigwedge_{q \in \mathbb{Q}_i} b_{\mathcal{M}}(I_{(i,j)}, q))}{|\mathbb{I}|} \quad (4)$$

Question Pair Accuracy:

$$qAcc = \frac{\sum_{i,k} \mathbb{1}(\bigwedge_{I \in \mathbb{I}_i} b_{\mathcal{M}}(I, q_{(i,k)}))}{|\mathbb{Q}|} \quad (5)$$

4.3. Analytical Evaluation Criteria

In addition to the accuracy metrics, we introduce three analytical criteria to measure and diagnose the failures of LVLMs, *Yes/No Bias Test*, *Consistency Test*, and *Diagnostic Test*. Instead of examining and analyzing each failed case qualitatively, we propose these novel quantitative measurements through the unique design of our question sets. These tests are listed in the order of complexity, so the latter test would not be as useful and insightful if the former basic test failed.

4.3.1 Yes / No Bias Test

According to [23], some models [15, 31, 50] tend to respond with “yes” in most cases. No further analysis is necessary if the model has a very strong bias or tendency to answer one way regardless of the actual question, so we design two criteria to reveal such preference of the model.

Yes Percentage Difference (Pct. Diff) $d_y \in [-1, 1]$:

$$d_y = \frac{\sum_{(I,q) \in \mathcal{V}} [\mathbb{1}(\mathcal{M}(I, q) = “yes”) - \mathbb{1}(y(I, q) = “yes”)]}{|\mathcal{V}|}, \quad (6)$$

d_y represents the difference between the predicted and actual number of “Yes” in the question set. The model is more biased when $|d_y|$ is close to 1.

False Positive Ratio (FP Ratio) $r_{fp} \in [0, 1]$:

$$r_{fp} = \frac{\sum_{(I,q) \in \mathcal{W}} \mathbb{1}(\mathcal{M}(I, q) = “yes”)}{|\mathcal{W}|}, \quad (7)$$

where $\mathcal{W} = \{(I, q) \in \mathcal{V} \mid b_{\mathcal{M}}(I, q) = 0\}$ is the set of incorrect visual questions. r_{fp} measures how likely the model responds with “Yes” out of all incorrect responses. The model is more robust when r_{fp} is close to 0.5.

4.3.2 Consistency Test

The goal of the consistency test is to test the logical consistency of responses and make sure questions are not answered based on random guesses. Many questions \mathbb{Q}^i from root \mathcal{R}^i are logically consistent: for example, “Is the left segment longer than/shorter than/equal to the right segment?” The consistency test is implemented and measured using *fAcc* (Metrics 4). We design the question set \mathbb{Q}_i to be logically correlated over a figure. Therefore, we consider the model *inconsistent* when only some of the questions in \mathbb{Q}_i are correct. In other cases, the model would be consistently correct or consistently wrong.

4.3.3 Language Hallucination and Visual Illusion

Before we dive into the diagnostic test, we categorize the failures into two major types based on the failed cases:

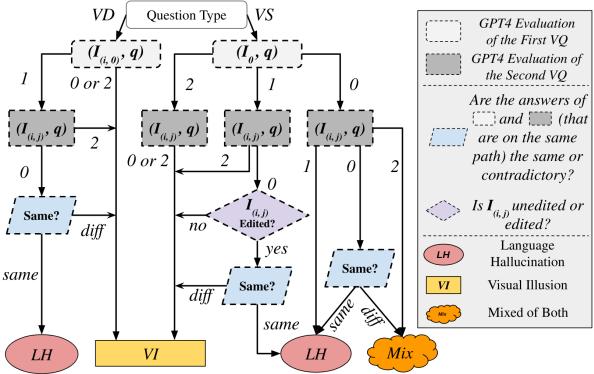


Figure 3. **Decision Tree to Diagnose Failure Types:** Based on the correctness of two questions in a control pair, and the difference of their responses, we use this decision tree to analyze the failure. The output of *GPT4 Evaluation* could be *Incorrect* (0), *Correct* (1), or *Uncertain* (2) if the predicted response is ambiguous.

Language Hallucination refers to perceptions formed without relevant visual input. In language hallucination, the model makes false prior assumptions about the input and image context based on its parametric memory. The model should respond based on how the question is framed instead of ignoring it or making false assumptions about the image.

Visual Illusion denotes the misinterpretation of accurate visual information. Visual illusion comes from the failure to recognize and understand the input image visually. The model could not obtain accurate information or reason about the image correctly.

4.3.4 Diagnostic Test

To study the issue of language hallucination and language illusion, we analyze the responses and correctness of both visual questions within a *VQ Control Pairs* and divide incorrect responses into three categories: *Language Hallucination*, *Visual Illusion*, and *Mixed / Uncertain*. We measure the percentage of those failures out of all failed cases.

Control Pair: The control pair will always contain an original image for *visual dependent* questions or an empty image (no visual) for *visual supplement* questions. The other question in the control pair may have an edited image (or an original image for VS question). The response to this question would provide more information on whether the answer exists in the parametric knowledge or if the model has seen it in the training data. In addition, we can examine whether the response remains the same after editing the original image to obtain more insights into the failures, which is more informative than checking a single visual question alone. In Fig. 3, we provide a decision tree to determine the type of failure for a control pair. We consider the following principles when assigning the failure types:

1. For *visual dependent* (VD) questions, or *visual supplement* (VS) questions that have visual inputs, if the re-

sponse is incorrect or uncertain, the failure could be **visual illusion**, since the model could not extract from the visual information correctly.

2. For *visual supplement* (VS) questions that don't have visual inputs, if the response gives a certain but wrong answer, we attribute it to **language hallucination**.
3. If the model responds to the original image (or no image) correctly and has the same response to the edited image (which is contrary to common sense), it means that the parametric knowledge overtakes the actual image input. Therefore, we also attribute the failure to **language hallucination**.

We will include some examples in the supplemental material.

5. Experimental Results

5.1. Models

We conduct massive experiments on HALLUSIONBENCH to evaluate a total of 15 LVLMs, including GPT-4V [1], LLaVA-1.5 [31], Gemini Pro Vision [39], Claude 3 [38], MiniGPT4 [63], MiniGPT5 [62], GiT [44], InstructBLIP [11], Qwen-VL [6], mPLUG-Owl-v1 [50], mPLUG-Owl-v2 [51], LRV-Instruction [28], BLIP2 [21], BLIP2-T5 [21], and Open-Flamingo [3]. We also include *Random Chance* (i.e. randomly choose Yes or No) as a baseline.

5.2. Result Analysis

We compare the performance of several models, including both closed-source models and open-sourced models. Results are given in Tab. 2, Tab. 3 and Fig. 4. Additionally, we established a human expert evaluation to assess the effectiveness of text-only GPT4-assisted evaluation.

Correctness Evaluation. As shown in Tab. 2, GPT-4V outperforms all the open-sourced LVLMs by a large margin except the *Hard Accuracy*. *Hard Accuracy* measures the models' ability to understand human-edited images from HALLUSIONBENCH. The poor accuracy demonstrates the challenges of our image manipulations for GPT-4V and other open-source LVLMs. In the open-sourced models, we investigate if expanding the size (0.8B to 13B) of the LLM backbone can mitigate object existence hallucination. As detailed in Tab. 2, there is a noticeable reduction in hallucination as the model size increases, like LLaVA-1.5 and BLIP2-T5. Among models with a size of less than 10B, InstructBLIP and mPLUG-Owl-v2 are the best-performing ones. InstructBLIP, leveraging the BLIP-2 architecture and enhanced through instruction fine-tuning across 26 diverse datasets, demonstrates that a broader and more extensive training set can substantially enhance performance. The boosting performance of mPLUG-Owl-v2 compared with mPLUG-Owl-v1 can be attributed to its novel module, which utilizes the language decoder acting as a universal interface for managing different modalities.

Method	# Parameter	Evaluation	Question Pair Accuracy (<i>qAcc</i>) \uparrow	Figure Accuracy (<i>fAcc</i>) \uparrow	Easy Accuracy (Easy <i>aAcc</i>) \uparrow	Hard Accuracy (Hard <i>aAcc</i>) \uparrow	All Accuracy (<i>aAcc</i>) \uparrow
GPT4V [1] (Oct 2023)	-	Human GPT4-Assisted	31.42 28.79	44.22 39.88	79.56 75.60	38.37 37.67	67.58 65.28
LLaVA-1.5 [31]	13B	Human GPT4-Assisted	9.45 10.55	25.43 24.86	50.77 49.67	29.07 29.77	47.12 46.94
Claude 3 [38]	-	GPT4-Assisted	21.76	28.61	55.16	41.40	56.86
Gemini Pro Vision [39] (Dec 2023)	-	GPT4-Assisted	7.69	8.67	35.60	30.23	36.85
BLIP2-T5 [21]	12.1B	GPT4-Assisted	15.16	20.52	45.49	43.49	48.09
Qwen-VL [6]	9.6B	GPT4-Assisted	5.93	6.65	31.43	24.88	39.15
Open-Flamingo [3]	9B	GPT4-Assisted	6.37	11.27	39.56	27.21	38.44
MiniGPT5 [62]	8.2B	GPT4-Assisted	10.55	9.83	36.04	28.37	40.30
MiniGPT4 [63]	8.2B	GPT4-Assisted	8.79	10.12	31.87	27.67	35.78
InstructBLIP [11]	8.2B	GPT4-Assisted	9.45	10.11	35.60	45.12	45.26
BLIP2 [21]	8.2B	GPT4-Assisted	5.05	12.43	33.85	40.70	40.48
mPLUG_Owl-v2 [51]	8.2B	GPT4-Assisted	13.85	19.94	44.84	39.07	47.30
mPLUG_Owl-v1 [50]	7.2B	GPT4-Assisted	9.45	10.40	39.34	29.77	43.93
LRV_Instruction [28]	7.2B	GPT4-Assisted	8.79	13.01	39.78	27.44	42.78
GIT [44]	0.8B	GPT4-Assisted	5.27	6.36	26.81	31.86	34.37
Random Chance	-	GPT4-Assisted	15.60	18.21	39.12	39.06	45.96

Table 2. **Correctness Leaderboard on HALLUSIONBENCH with various LVLMs:** All the numbers are presented in % and the full score is 100%. Hard questions refer to the edited images. We highlight the Top 3 models with the GPT4-assisted evaluation.

Method	# Parameter	Evaluation	Yes/No Bias		Consistency			Language and Vision Diagnosis		
			Pct. Diff (~ 0)	FP Ratio (~ 0.5)	Correct \uparrow	Inconsistent \downarrow	Wrong \uparrow	Language Hallucination	Visual Illusion	Mixed
GPT4V [1] (Oct 2023)	-	Human GPT4-Assisted	0.066 0.058	0.60 0.58	44.22 39.88	32.66 38.15	23.12 21.97	21.86 22.19	46.17 45.66	31.97 32.14
LLaVA-1.5 [31]	13B	Human GPT4-Assisted	0.27 0.26	0.76 0.75	25.43 24.86	42.49 45.38	32.08 29.77	25.63 26.71	51.42 51.09	22.95 22.20
Claude 3 [38]	-	GPT4-Assisted	0.063	0.57	28.61	49.42	21.97	19.10	59.14	21.77
Gemini Pro Vision [39] (Dec 2023)	-	GPT4-Assisted	-0.02	0.48	8.67	56.94	34.39	25.95	49.37	24.68
BLIP2-T5 [21]	12.1B	GPT4-Assisted	0.08	0.58	20.52	59.54	19.94	41.64	40.44	17.92
Qwen-VL [6]	9.6B	GPT4-Assisted	0.12	0.60	6.65	50.29	43.06	0.87	88.06	11.06
Open-Flamingo [3]	9B	GPT4-Assisted	0.33	0.77	11.27	59.83	28.90	30.07	48.06	21.87
MiniGPT5 [62]	8.2B	GPT4-Assisted	0.28	0.71	9.83	56.36	33.82	10.09	73.44	16.47
MiniGPT4 [63]	8.2B	GPT4-Assisted	0.19	0.65	10.12	57.80	32.08	23.59	56.55	19.86
InstructBLIP [11]	8.2B	GPT4-Assisted	-0.13	0.38	10.12	68.50	21.39	29.29	54.53	16.18
BLIP2 [21]	8.2B	GPT4-Assisted	0.18	0.65	12.43	63.01	24.57	39.14	43.45	17.41
mPLUG_Owl-v2 [51]	8.2B	GPT4-Assisted	0.25	0.77	19.94	58.09	21.97	28.24	50.42	21.34
mPLUG_Owl-v1 [50]	7.2B	GPT4-Assisted	0.32	0.79	10.40	60.12	29.48	3.95	78.36	17.69
LRV_Instruction [28]	7.2B	GPT4-Assisted	0.26	0.73	13.01	53.47	33.53	4.49	76.47	19.04
GIT [44]	0.8B	GPT4-Assisted	0.04	0.53	6.36	53.76	39.88	30.90	58.30	10.80
Random Chance	-	GPT4-Assisted	0.08	0.57	18.20	57.51	24.28	-	-	-

Table 3. **Analytical Evaluation Results on HALLUSIONBENCH with various LVLMs:** *Pct. Diff* ranges from [-1, 1]. The model is more biased when *Pct. Diff* is close to -1 or 1. *FP Ratio* ranges from [0, 1]. The model is more robust when *FP Ratio* is close to 0.5. All the other metrics are presented in %, and the full score is 100%. We highlight the Top 3 models with the GPT4-assisted evaluation.

Yes/No Bias. Another observation is that GPT-4V, BLIP2-T5, and mPLUG-Owl-v2 outperform *Random Choice* in both question pair accuracy, figure pair accuracy, and question level accuracy. Other models, such as Qwen-VL and MiniGPT4, perform even worse than *Random Choice*. This indicates their visual reasoning abilities are still limited. However, LLaVA-1.5 outperforms *Random Choice* while achieving poor results in both question pair accuracy and figure pair accuracy. We attribute this phenomenon to the fact that LLaVA-1.5 tends to answer *Yes*. This assumption is supported by the low *Yes Percentage Difference* and *False Positive Ratio* of LLaVA-1.5 in *Yes/No Bias Test* from Tab. 3. Besides, we find that Open-Flamingo and mPLUG-Owl-v1 also tend to answer *Yes* with the high *Yes Percentage Differ-*

ence and *False Positive Ratio*. Inspired by [28], one possible reason is that these LVLMs lack balanced positive and negative instructions in their training set. We also attribute the poor performance of these LVLMs to the scarcity of human-edited images in their training set since most LVLMs only utilize original images from existing datasets.

Language and Vision Diagnosis. We report fine-grained scores of six prominent LVLMs across different visual inputs in Fig. 4. Results show that *Math*, *Illusion*, and *Video* is the most challenging format for current LVLMs, including GPT-4V. From Fig. 5 (top), we found both GPT-4V and LLaVA-1.5 are unable to correctly recognize regular triangles, meaning that geometry and math are still a challenging task for GPT-4V. From Fig. 5 (middle), we found GPT-4V is

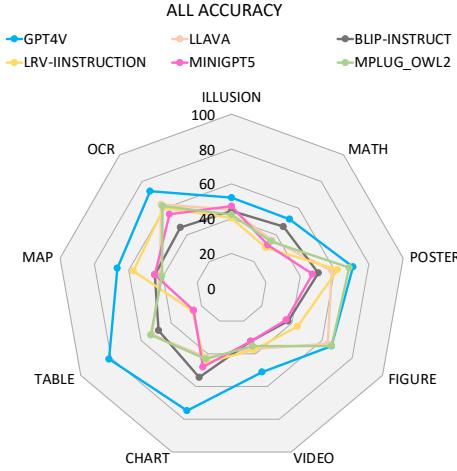


Figure 4. **Accuracies on each subcategories:** We show six prominent LVLMs on HALLUSIONBENCH across different types.

more knowledgeable than LLaVA-1.5 in recognizing all the illusion cases and knowing their names. However, GPT-4V fails to answer the question faithfully based on the edited images. The reason behind this might be that GPT-4V tends to generate answers based on its parametric memory instead of analyzing the images. Compared to GPT-4V, LLaVA-1.5 performs badly on both the original image and edited images, indicating that the visual perception skill of LLaVA-1.5 is limited. From Fig. 5 (bottom), we found that GPT-4V is unable to distinguish between the positive sequence and the reversed sequence of the images, indicating that there is still much room to improve the video reasoning ability.

6. Conclusion, Limitations and Future Work

In this work, we introduce HALLUSIONBENCH, the first advanced diagnostic suite to analyze the failure cases of 15 current LVLMs. HALLUSIONBENCH presents significant challenges to existing LVLMs like GPT-4V(ision), by emphasizing nuanced understanding and interpretation of visual data. Moreover, our unique design of the visual-question pairs facilitates a quantitative analysis of the models’ failures, enabling a more thorough evaluation. We share our observations and key insights for future studies:

- When GPT-4V, LLaVA-1.5, and other LVLMs have prior knowledge of questions in HALLUSIONBENCH, they usually suffer from Language Hallucination as they tend to prioritize their prior knowledge which leads to incorrect answers. The model should handle the trade-off between parametric memory and context.
- When LVLMs have not had parametric memory or prior knowledge regarding the questions in HALLUSIONBENCH, they can still be prone to Visual Illusion and prefer to produce wrong answers about the given figure. The visual capability of existing LVLMs is still limited.
- GPT-4V and other LVLMs can be easily misled by simple image manipulations in HALLUSIONBENCH, including

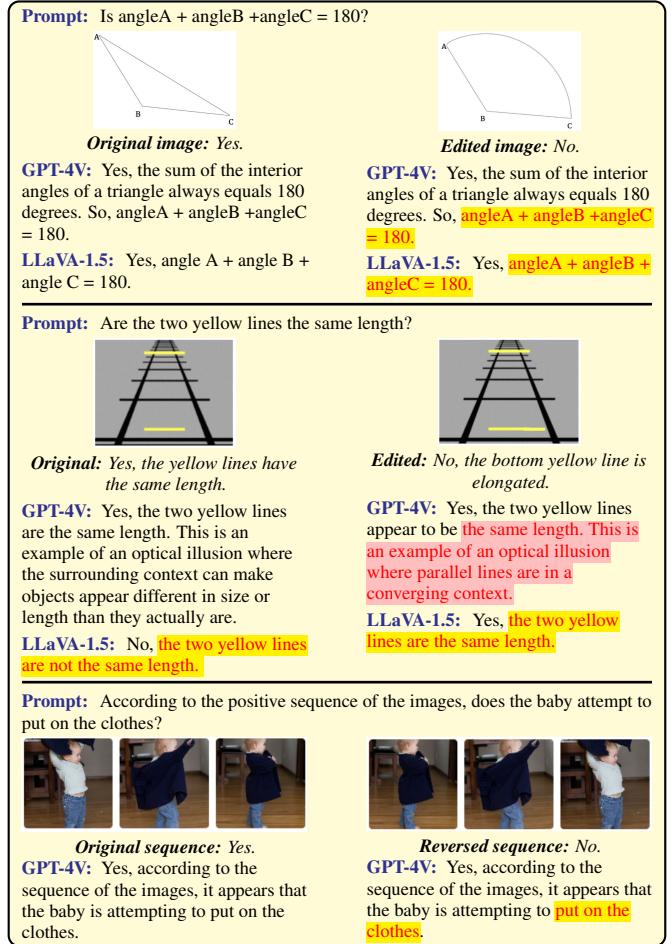


Figure 5. **Failure Cases in Math, Illusion and Video:** We highlight language hallucination and visual illusion.

image flipping, order reversing, masking, optical character editing, object editing, and color editing.

- GPT-4V and other LVLMs are unable to capture the temporal relations of multiple images and fail to answer temporal reasoning questions in HALLUSIONBENCH. The existing LVLMs lack true temporal reasoning ability.

We plan to expand this benchmark and figure out other ways to diagnose issues within LVLMs. We hope that HALLUSIONBENCH can be used to identify and provide insights on the weakness of different LVLMs, to facilitate finetuning and improvement of those models based on the diagnoses.

7. Acknowledgements

This research was supported by Army Cooperative Agreement W911NF2120076 and ARO W911NF2310046 and W911NF2310352. Our work is also supported in part by DARPA SemaFor Program under HR001120C0124. Zhou is supported in part by Adobe Research gift fund. Xiaoyu and Huang are supported by NSF-IIS-2147276 FAI, DOD N00014-22-1-2335 and FA9550-23-1-0048, DARPA GARD HR00112020007, Adobe, Capital One and JP Morgan.

A. More Case Analysis on HALLUSIONBENCH with GPT-4V and LLaVA-1.5

In this section, we give a few samples in HALLUSIONBENCH and share our observations. **Each figure is self-contained for readability**, where we highlight the control pairs, the responses of GPT-4V and LLaVA-1.5, the failures of those models, and the corresponding part of the answers.

A.1. Visual Dependent Examples

From the famous illusions in Fig.7, Fig.8, and Fig.9, we found GPT-4V is more knowledgeable than LLaVA-1.5 in recognizing all the illusion cases and knowing their names. However, GPT-4V fails to answer the question faithfully based on the edited images. The reason behind this might be that GPT-4V tends to generate answers based on its parametric memory instead of analyzing the images. Compared to GPT-4V, LLaVA-1.5 performs badly on both the original image and edited images, indicating that the visual perception skill of LLaVA-1.5 is limited.

From the examples in Fig.10 and Fig.11, we found both GPT-4V and LLaVA-1.5 are unable to correctly recognize parallel lines, regular triangles, polygons, and other math theorems, meaning that geometry and math are still a challenging task for GPT-4V.

We further explore GPT-4V's and LLaVA-1.5's abilities in Optical Character Recognition in Fig.12 and Figure Recognition in Fig.13. From our observations, we found that GPT-4V and LLaVA-1.5 are easily misled by editing the characters in the images, demonstrating that GPT-4V and LLaVA-1.5 generate answers based on their parametric memory instead of visual reasoning. This is because the difference between the original images and edited images is obvious.

Inspired by [48], which shows the promising video understanding of GPT-4V, we also investigate more examples in Fig.14 and Fig.15, including several frame sequence examples. The positive sequence and reversed sequence have the opposite semantic meaning, such as "*disappear or appear*" and "*park or leave*" in Fig.14. From the comparison, we found that GPT-4V is unable to distinguish between the positive sequence and the reversed sequence of the images, indicating that there is still much room to improve the video reasoning ability.

A.2. Visual Supplement Examples

In Fig.16, Fig.17, and Fig.18, GPT-4V does not have an affirmative answer if no images are given. Given the image context, GPT-4V and LLaVA-1.5 are unable to understand the chart correctly, indicating that their chart reasoning ability is still limited. In the second example (bottom) of Fig.24, the predictions of GPT-4V changed completely after we rotated the chart.

In Fig.19, Fig.20, Fig.22, Fig.23, and Fig.24, GPT-4V and LLaVA-1.5 have an affirmative answer if no images are given. After providing the image, including charts, tables, or maps, we found that they preferred to answer the questions with their knowledge instead of analyzing the image. This might be because GPT-4V and LLaVA-1.5 demonstrate a marked dependence on textual reasoning capabilities, often prioritizing them over visual reasoning.

From Fig. 20 and Fig.21, we found the knowledge from LLaVA-1.5 is not accurate since it states " π doesn't range from 3.1415926 and 3.1415927" and "North Carolina is farther north than Delaware." This observation also supports our claim that GPT-4V is more knowledgeable than LLaVA-1.5.

B. Decision Tree Logic and Examples

In Fig. 6, we utilize the decision tree to determine the failure types. In the rest of the section, specifically Fig. 25-36, we will provide a few examples and explain the logic that leads to different types of errors. **Each figure with its caption is self-contained for readability**.

In Fig. 25 (bottom), it is a visual-dependent sample (VD). The answer regarding the original image is correct (1), but the answer to the edited image is incorrect (0), and the two answers are the same (*same*). This shows that GPT-4V knows the "*Chubb illusion*" in its parametric knowledge but can not answer according to the image. In Fig. 6, these correspond to the (VD) R-G-R-C route in the decision tree, leading to the diagnostic result of *Language Hallucination*.

In Fig. 26 (bottom), it is a visual-dependent sample (VD). The answer regarding the original image is correct (1), but the answer to the edited image is incorrect (0), and the two answers are not the same (*same*). This shows that GPT-4V can not compare the length of the two lines correctly. In Fig. 6, it corresponds to the (VD) R-G-R-M-B route in the decision tree, leading to the diagnostic result of *Visual Illusion*.

In Fig. 27 (bottom), it is a visual-dependent sample (VD). The answer regarding the original image is correct (1), but the answer to the edited image is uncertain (2). This shows that GPT-4V is uncertain about the length of the vertical line compared with the horizontal line. In Fig. 6, it corresponds to the (VD) R-G-B-B route in the decision tree, leading to the diagnostic result of *Visual Illusion*.

In Fig. 28 (bottom), It is a visual-dependent sample (VD). The answer regarding the original image is incorrect (0) or uncertain (2). This shows that LLaVA-1.5 fails to determine the diameters of the three circles in the original image, but succeeds in the edited image. In Fig. 6, it corresponds to the (VS) R-B route in the decision tree, leading to the diagnostic result of *Visual Illusion*.

In Fig. 29 (bottom), it is a visual-supplement sample (VS). The answer regarding the original image is uncertain

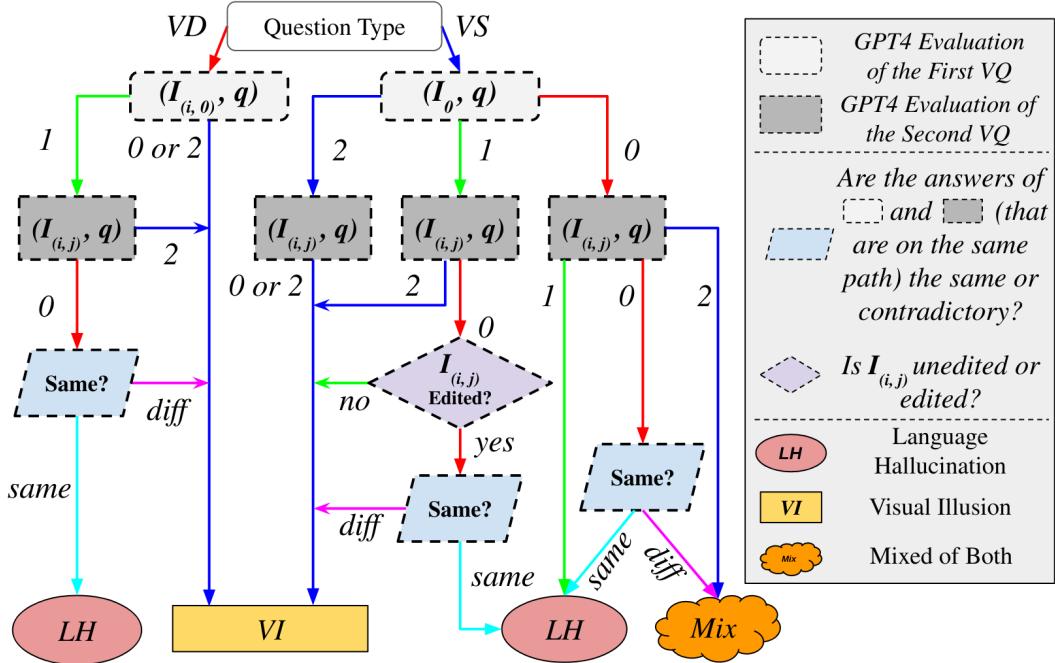


Figure 6. **Decision Tree to Diagnose Failure Types:** Based on the correctness of two questions in a control pair, and the difference in their responses, we use this decision tree to analyze the failure. We highlight different decision paths with Red(R), Blue(B), Green(G), Cyan(C) and Magenta(M). So a path on the decision tree can be represented as a sequence of colors, e.g., R-G-R-C. The output of *GPT4 Evaluation* could be *Incorrect* (0), *Correct* (1), or *Uncertain* (2) if the predicted response is ambiguous.

(2), but the answer is incorrect (0) or uncertain (2) when the supplementary image is given. This shows that GPT-4V is uncertain about the answer without the visual input, and fails to answer the question with the supplementary image as well. In Fig. 6, it corresponds to the (VS) B-B-B route in the decision tree, leading to the diagnostic result of *Visual Illusion*.

In Fig. 30 (bottom), It is a visual-supplement sample (VS). The answer is correct (1) without being given any image. However, the answer is uncertain (2) when the supplementary image is given. This shows that GPT-4V is uncertain about the answer given the supplementary image though it could make the correct answer without the image. In Fig. 6, it corresponds to the (VS) B-G-B-B route in the decision tree, leading to the diagnostic result of *Visual Illusion*.

In Fig. 31 (bottom), it is a visual-supplement sample (VS). The answer is already correct (1) without being given any image. However, the answer is incorrect (0) given the original supplementary image. The supplementary image is not edited. This shows that GPT-4V produces the wrong answer given the supplementary image, though it could produce the correct answer without the image. In Fig. 6, it corresponds to the (VS) B-G-R-G-B route in the decision tree, leading to the diagnostic result of *Visual Illusion*.

In Fig. 32 (bottom), it is a visual-supplement sample

(VS). The answer is correct (1) without being given any image. However, the answer is incorrect (0) when a edited image is given. The supplementary image is edited and the two answers are not the same. This shows that GPT-4V produces the wrong answer based on reasons inconsistent with the edited supplementary image, though it could produce a correct answer without the image. In Fig. 6, it corresponds to the (VS) B-G-R-R-M-B route in the decision tree, leading to the diagnostic result of *Visual Illusion*.

In Fig. 33 (bottom), it is a visual-supplement sample (VS). The answer is correct (1) without being given any image but the answer is incorrect (0) when an edited supplementary image is given. The supplementary image is edited by swapping Delaware and Arizona on the map. The two answers are the same. This indicates that GPT-4V has the prior knowledge of “Delaware is the farthest north” in its parametric knowledge but can not provide a correct answer according to the edited map. In Fig. 6, it corresponds to the (VS) B-G-R-R-C route in the decision tree, leading to the diagnostic result of *Language Hallucination*.

In Fig. 34 (bottom), it is a visual-supplement sample (VS). The answer is incorrect (0) without being given any image. But the answer becomes correct given the original image. This indicates that LLaVA-1.5’s answer is affected by hallucinations without given image information. In Fig. 6,

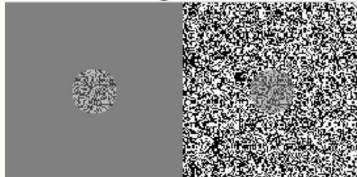
it corresponds to the (VS) **B-R-G** route in the decision tree, leading to the diagnostic result of *Language Hallucination*.

In Fig. 35 (bottom), it is a visual-supplement sample (VS). The answer is incorrect (0) without being given any image. The answer is still incorrect (0) when the original supplementary image is given. And the two answers are the same. This shows that LLaVA-1.5 has the issue of hallucinations with and without the image information. In Fig. 6, it corresponds to the (VS) **B-R-R-C** route in the decision tree, leading to the diagnostic result of *Language Hallucination*.

In Fig. 36 (bottom), it is a visual-supplement sample (VS). The answer is incorrect (0) without being given any image. The answer is still incorrect (0) when an edited supplementary image is given. However, the two answers are not the same. This indicates that the commonsense knowledge about the location of US states in LLaVA-1.5 is weak and wrong without the input image of the US map. Additionally, the visual interpretation of the map by LLaVA-1.5 is incorrect. In Fig. 6, it corresponds to the (VS) **B-R-R-M** route in the decision tree, leading to the diagnostic result of *Potentially Mixed*.

Visual Dependent – Optical Illusion of Size

Prompt: The image shows the famous Chubb illusion. Are the center areas of two rectangular fields identical?

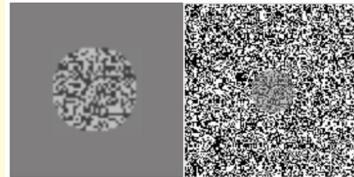


Original image

Answer: Yes.

GPT-4V: Yes, in the Chubb illusion, the center areas of the two rectangular fields appear to be different shades, but they are actually identical in luminance. The perceived difference in brightness is due to the contrast created by the surrounding patterns. This illusion highlights how context can influence our perception.

LLaVA-1.5: Yes, the center areas of the two rectangular fields are identical.



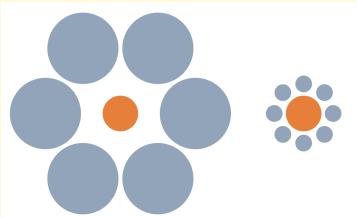
Edited: The circle on the left are enlarged.

Answer: No, they are not identical.

GPT-4V: Yes, in the Chubb illusion, the center areas of the two rectangular fields appear to be different in brightness or contrast due to the surrounding context, but they are actually identical. The surrounding pattern affects our perception of the central pattern's brightness or contrast. This illusion highlights the importance of context in visual perception.

LLaVA-1.5: Yes, the center areas of two rectangular fields are identical.

Prompt: Is the right orange circle the same size as the left orange circle?

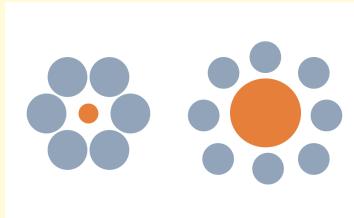


Original

Answer: Yes, the orange balls have the same size.

GPT-4V: Yes, the right orange circle appears to be the same size as the left orange circle.

LLaVA-1.5: No, the right orange circle is smaller than the left orange circle.



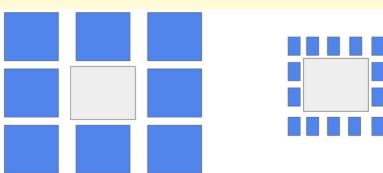
Edited: The orange ball on the right is enlarged.

Answer: No, the orange balls have different size.

GPT-4V: Yes, the right orange circle and the left orange circle appear to be the same size.

LLaVA-1.5: Yes, the right orange circle is the same size as the left orange circle.

Prompt: Is the right grey box larger than the left grey box?

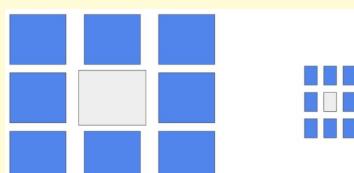


Original

Answer: No, the grey boxes have the same size.

GPT-4V: No, both grey boxes appear to be of the same size.

LLaVA-1.5: Yes, the right grey box is larger than the left grey box.



Edited

Answer: No, the right grey box is smaller than the left.

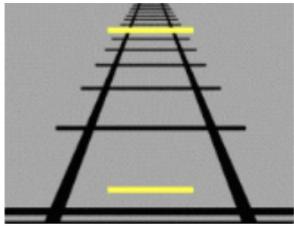
GPT-4V: No, the right grey box and the left grey box appear to be the same size.

LLaVA-1.5: Yes, the right grey box is larger than the left grey box.

Figure 7. We highlight the incorrect answer according to visual illusion, language hallucination, or potentially mixed. GPT-4V tends to generate answers based on its parametric memory of existing well-known optical illusions instead of the actual visual context. Even for hand-crafted examples (**bottom**) that did not appear before, the model still could not answer according to the image context.

Visual Dependent – Optical Illusion of Length

Prompt: Are the two yellow lines the same length?

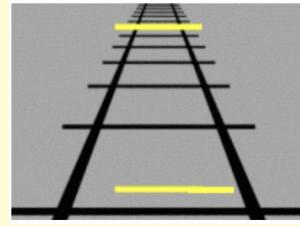


Original

Answer: Yes, the yellow lines have the same length.

GPT-4V: Yes, the two yellow lines are the same length. This is an example of an optical illusion where the surrounding context can make objects appear different in size or length than they actually are.

LLaVA-1.5: No, **the two yellow lines are not the same length.**



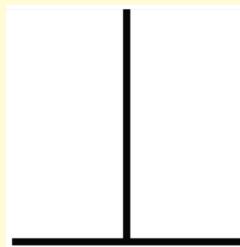
Edited: The bottom yellow line is elongated.

Answer: No, the bottom yellow line is longer.

GPT-4V: Yes, the two yellow lines appear to be **the same length. This is an example of an optical illusion where parallel lines are in a converging context.**

LLaVA-1.5: Yes, **the two yellow lines are the same length.**

Prompt: Is the vertical line the same length as the horizontal line?

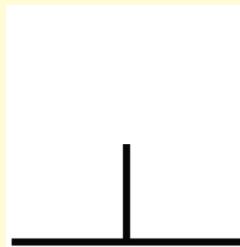


Original

Answer: Yes, the lines have the same length.

GPT-4V: Yes, the vertical line is the same length as the horizontal line.

LLaVA-1.5: Yes, the vertical line is the same length as the horizontal line.



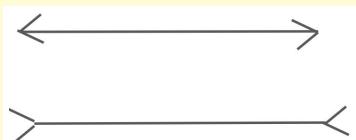
Edited: The vertical line is shorten.

Answer: No, the vertical line is shorter.

GPT-4V: **No, the vertical line appears to be longer than the horizontal line in the image provided.**

LLaVA-1.5: Yes, **the vertical line is the same length as the horizontal line.**

Prompt: Is the top line longer than the bottom line?

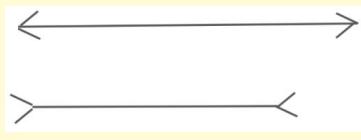


Original

Answer: No, the lines have the same length.

GPT-4V: No, both lines appear to be of the same length. The arrowheads at the ends can create an optical illusion that might make them seem different in length, but the lines are the same.

LLaVA-1.5: Yes, **the top line is longer than the bottom line.**



Edited: The top line is elongated.

Answer: Yes, the top line is longer.

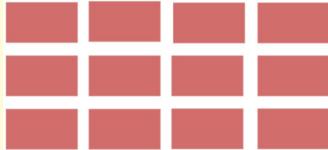
GPT-4V: Both lines appear to be of **the same length.** The arrowhead style at the ends might create an **optical illusion,** but the lines themselves seem equal in length.

LLaVA-1.5: Yes, the top line is longer than the bottom line.

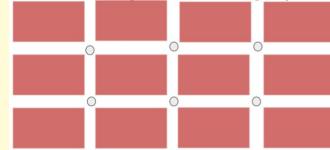
Figure 8. We highlight the incorrect answer according to **visual illusion**, **language hallucination**, or **potentially mixed**. GPT-4V can recognize many optical illusion cases but is also easily tricked by the scene and setup of the images. Both models have bad performance in recognizing and measuring length.

Visual Dependent – Other Illusion

Prompt: Are there any grey points actually in the intersection of the white lines?



Original image



Edited: Grey points are added.

Answer: No.

GPT-4V: No, there are no grey points in the image you provided. The image contains only red rectangles arranged in a grid pattern.

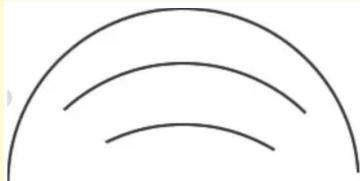
LLaVA-1.5: Yes, **there are grey points in the intersection of the white lines.**

Answer: Yes.

GPT-4V: No, **there are no grey dots in the intersections.** The optical illusion you're referencing is known as the **Hermann grid illusion**. When looking at the grid, people often perceive ghost-like grey blobs at the intersections, but when they try to look directly at them, the blobs disappear.

LLaVA-1.5: Yes, there are grey points in the intersection of the white lines. These grey points are located at the corners of the red squares.

Prompt: According to the Torsion dystonia illusion, are the diameters of the three circles the same?

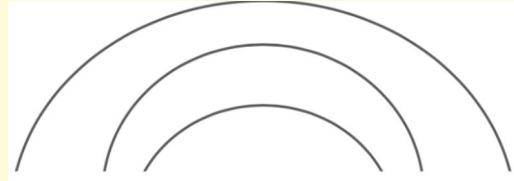


Original

Answer: Yes, the curves have the same diameters.

GPT-4V: Yes, according to the Torsion Dystonia illusion, the diameters of the three circles are actually the same. The illusion tricks our perception, making some circles appear larger or smaller than they really are.

LLaVA-1.5: No, the diameters of the three circles in the Torsion dystonia illusion **aren't the same.**



Edited: The curves have different diameters.

Answer: No, the curves have different diameters.

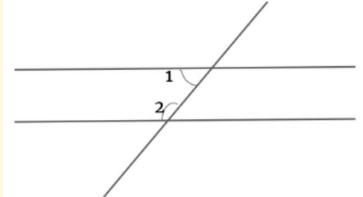
GPT-4V: Yes, **all three arcs are part of circles with the same diameter.** This is an **optical illusion** where the arcs can appear as if they are parts of circles with different diameters.

LLaVA-1.5: No, the diameters of the three circles in the Torsion dystonia illusion are not the same. The illusion is designed to make it appear as if the circles are of the same size, but they are actually different in size.

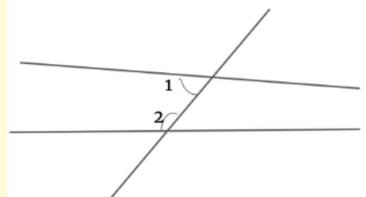
Figure 9. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. GPT-4V recognizes the illusion cases but fails to answer the question faithfully based on the actual image context.

Visual Dependent – Geometry&Math

Prompt: According to parallel lines theorem, is angle 1 + angle 2 = 180 ?



Original image



Edited image

Answer: Yes.

GPT-4V: Yes, according to the properties of parallel lines cut by a transversal, angle 1 and angle 2 are supplementary. Therefore, the sum of angle 1 and angle 2 is 180 degrees.

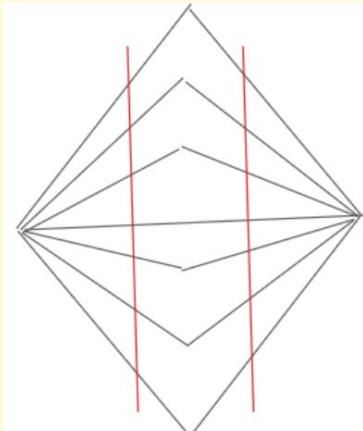
LLaVA-1.5: Yes, according to the parallel lines theorem, the sum of the angles formed by two parallel lines is always 180 degrees. In the image, the two parallel lines are labeled with angles 1 and 2, and their sum is indeed 180 degrees.

Answer: No.

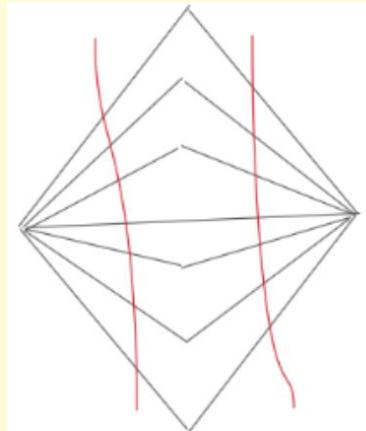
GPT-4V: Yes, according to the properties of parallel lines cut by a transversal, angle 1 and angle 2 are supplementary. This means that the **sum of angle 1 and angle 2 is 180 degrees**.

LLaVA-1.5: Yes, according to the parallel lines theorem, the sum of the angles formed by two parallel lines is always 180 degrees. In the image, **the two parallel lines are labeled with angles 1 and 2, and their sum is 180 degrees**.

Prompt: Are the two red vertical lines both straight?



Original image



Edited image

Answer: Yes.

GPT-4V: Yes, the two red vertical lines are both straight.

LLaVA-1.5: Yes, the two red vertical lines are both straight.

Answer: No.

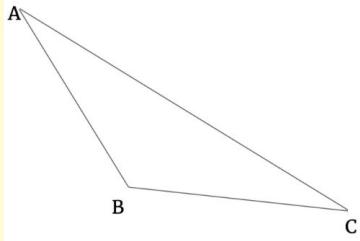
GPT-4V: Yes, the **two red vertical lines are both straight**.

LLaVA-1.5: No, one of the red vertical lines is curved, while the other one is straight.

Figure 10. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. **Top:** GPT-4V and LLaVA-1.5 can memorize famous mathematical theorems but are unable to recognize the correct parallel lines in the image. **Bottom:** GPT-4V is unable to distinguish whether two lines are straight. We attribute this failure to the lack of geometry recognition ability.

Visual Dependent – Geometry&Math

Prompt: Is angleA + angleB +angleC = 180?

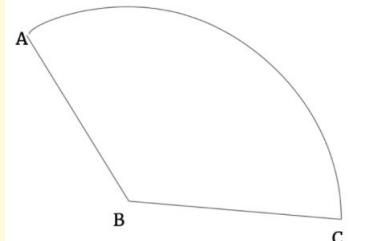


Original image

Answer: Yes.

GPT-4V: Yes, the sum of the interior angles of a triangle always equals 180 degrees. So, angleA + angleB +angleC = 180.

LLaVA-1.5: Yes, angle A + angle B + angle C = 180.



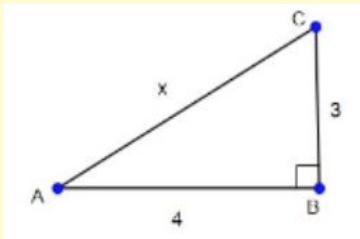
Edited image

Answer: No.

GPT-4V: Yes, the sum of the interior angles of a triangle always equals 180 degrees. So, **angleA + angleB +angleC = 180.**

LLaVA-1.5: Yes, **angleA + angleB + angleC = 180.**

Prompt: Is the value of x in the image 5?

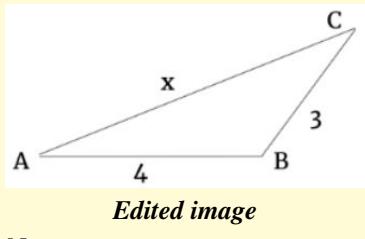


Original image

Answer: Yes.

GPT-4V: Yes, the value of x is 5. Given the triangle ABC in the image, it's a right triangle (as indicated by the square at angle C). Using the Pythagorean theorem, $x=5$.

LLaVA-1.5: Yes, the value of x in the image is 5.



Edited image

Answer: No.

GPT-4V: To determine the value of x in the triangle, we can use the Pythagorean theorem for right triangles, Where a and b are the lengths of the two shorter sides (legs) and c is the length of the hypotenuse. So, yes, the **value of x in the image is 5.**

LLaVA-1.5: Yes, **the value of x in the image is 5.**

Figure 11. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. In these examples, we modify important geometric properties of the triangles, and neither GPT-4V nor LLaVA-1.5 can recognize those changes. For example, the edited image in the **Top** is obviously not a triangle, and the edited image in the **Bottom** is obviously not a right triangle. We attribute this failure to the lack of geometry recognition ability.

Visual Dependent – Optical Character Recognition

Prompt: Does the image show "Beijing Roast Duck"?



Original image

Answer: Yes.

GPT-4V: Yes, the image appears to depict "Beijing Roast Duck" (also known as Peking Duck).

LLaVA-1.5: Yes, the image shows a picture of "Beijing Roast Duck," which is a popular dish in Chinese cuisine.



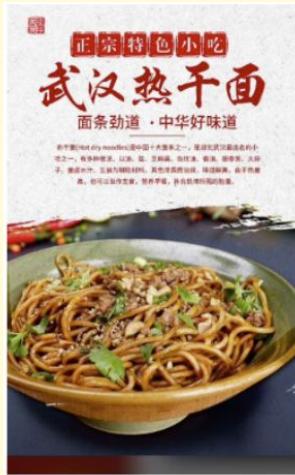
Edited image

Answer: No.

GPT-4V: Yes, the image appears to depict "Beijing Roast Duck", a famous Chinese dish known for its crispy skin and tender meat. The Chinese characters at the top of the image also confirm this as they read "Beijing Roast Duck."

LLaVA-1.5: Yes, the image shows a picture of "Beijing Roast Duck," which is a popular dish in Chinese cuisine.

Prompt: Does the image show "Hot dry noodles from Wuhan"?



Original image

Answer: Yes.

GPT-4V: Yes, the image displays a dish of noodles with the Chinese text indicating "Hot dry noodles", which is a traditional dish from Wuhan.

LLaVA-1.5: Yes, the image shows a bowl of "Hot dry noodles from Wuhan" with Chinese characters on the side.



Edited image

Answer: No.

GPT-4V: Yes, the image shows "Hot dry noodles - the taste of Wuhan." So, the image does indeed depict "Hot dry noodles from Wuhan."

LLaVA-1.5: Yes, the image shows a bowl of hot dry noodles from Wuhan.

Figure 12. We highlight the incorrect answer according to visual illusion, language hallucination, or potentially mixed. We highlight several advertisements with famous regional dishes with modifications on the regions. In both cases, GPT-4V and LLaVA-1.5 ignore the context and still reply with the well-known regions for that food.

Visual Dependent – Figure Recognition

Prompt: In the image, is Xuanzang accompanied by three disciples?



Original image

Answer: Yes.

GPT-4V: Yes, in the image, Xuanzang appears to be accompanied by his three disciples: Sun Wukong (the Monkey King), Zhu Bajie (the Pig), and Sha Wujing (the Sand Monk). This depiction is consistent with the classic Chinese novel "Journey to the West".

LLaVA-1.5: Yes, in the image, Xuanzang is accompanied by three disciples, who are depicted as a group of men standing next to each other.



Edited: One of the disciples is replaced.

Answer: No.

GPT-4V: Yes, in the image, Xuanzang appears to be accompanied by his three disciples: Sun Wukong (the Monkey King), Zhu Bajie (the Pig), and Sha Wujing (the Sand Monk). This depiction is consistent with the classic Chinese novel "Journey to the West".

LLaVA-1.5: Yes, in the image, Xuanzang is accompanied by three disciples, including a girl and two men.

Prompt: Is it the logo for Air Jordan?



Air Jordan Logo

Answer: Yes.

GPT-4V: Yes, that is the logo for Air Jordan, which is associated with Michael Jordan and the Nike brand.

LLaVA-1.5: Yes, it is the logo for Air Jordan, which is a popular sports brand associated with basketball and athletic footwear.



QIAODAN Logo (from China)

Answer: No.

GPT-4V: Yes, the image you provided resembles the Air Jordan logo, which is associated with Michael Jordan and the brand of basketball shoes and athletic clothing produced by Nike.

LLaVA-1.5: Yes, the image is the logo for Air Jordan, which is a sports brand known for its basketball shoes and apparel.

Figure 13. We highlight the incorrect answer according to visual illusion, language hallucination, or potentially mixed. **Top:** The judgments of GPT-4V and LLaVA-1.5 are affected by parametric memory and stereotyped judgment, meaning that they do not recognize that the girl is not Xuanzang's disciple. LLaVA-1.5 detects a girl and two men in the image but still regards her as the disciple of Xuanzang. **Bottom:** GPT-4V uses "resembles" due to the uncertainty of the answer but still mistakenly associates it with the Air Jordan brand. We attribute this failure to the Language Hallucination.

Visual Dependent – Video/Multiple Images

Prompt: According to the positive sequence images, does Homer Simpson disappear into the bushes?



Original sequence

Answer: Yes.

GPT-4V: Yes, the sequence of images shows Homer Simpson gradually disappearing into the bushes.

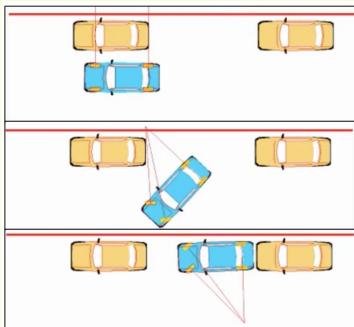


Reversed sequence

Answer: No.

GPT-4V: Yes, the sequence of images shows Homer Simpson gradually disappears into the bushes.

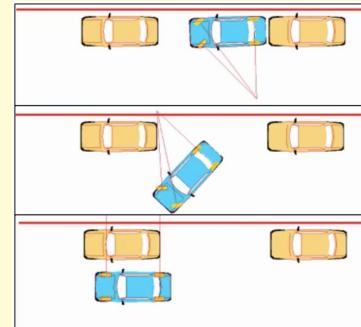
Prompt: According to the positive sequence of the images, does the blue car attempt to park?



Original sequence (From Top to Bottom)

Answer: Yes.

GPT-4V: Yes, according to the sequence of the images, the blue car appears to be attempting to parallel park between the two yellow cars.



Reversed sequence (From Top to Bottom)

Answer: No.

GPT-4V: Yes, according to the sequence of the images, the blue car appears to be attempting to parallel park between the two yellow cars.

Figure 14. We highlight the incorrect answer according to **visual illusion**, **language hallucination**, or **potentially mixed**. From the examples, we found that GPT-4V lacks the ability of temporal reasoning and is unable to distinguish between the positive sequence and reversed sequence of the images, even though they represent different actions.

Visual Dependent – Video/Multiple Images

Prompt: According to the positive sequence of the images, is the plug put in the power outlet?



Original sequence

Answer: Yes.

GPT-4V: Yes, the plug is put in the power outlet.



Reversed sequence

Answer: No.

GPT-4V: In the positive sequence of the images, the plug is **put into the power outlet**.

Prompt: According to the positive sequence of the images, does the baby attempt to put on the clothes?



Original sequence

Answer: Yes.

GPT-4V: Yes, according to the sequence of the images, it appears that the baby is attempting to put on the clothes.

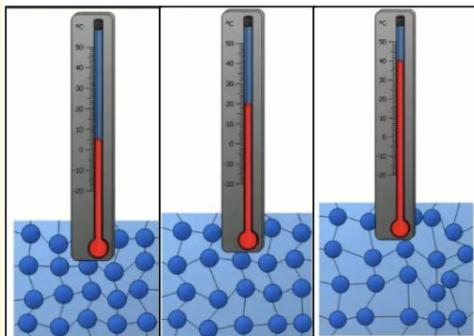


Reversed sequence

Answer: No.

GPT-4V: Yes, according to the sequence of the images, it appears that the baby is attempting to **put on the clothes**.

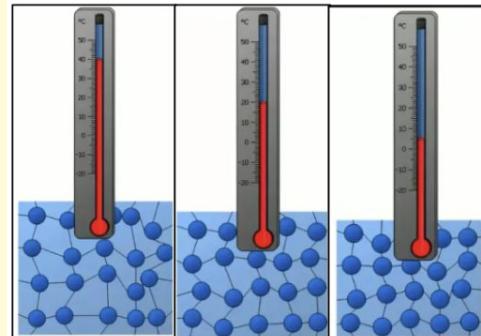
Prompt: According to the positive sequence of the images, does the temperature increase?



Original sequence

Answer: Yes.

GPT-4V: Yes, according to the sequence of the images, the temperature increases.



Reversed sequence

Answer: No.

GPT-4V: Yes, according to the positive sequence of the images, **the temperature increases**.

Figure 15. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. From the examples, we found that GPT-4V lacks the ability of temporal reasoning and is unable to distinguish between the positive sequence and reversed sequence of the images, even though they represent different actions.

Visual Supplement – Chart

Prompt: Is China, Hongkong SAR, the leading importing country of gold, silverware, and jewelry with the highest import value in 2018?

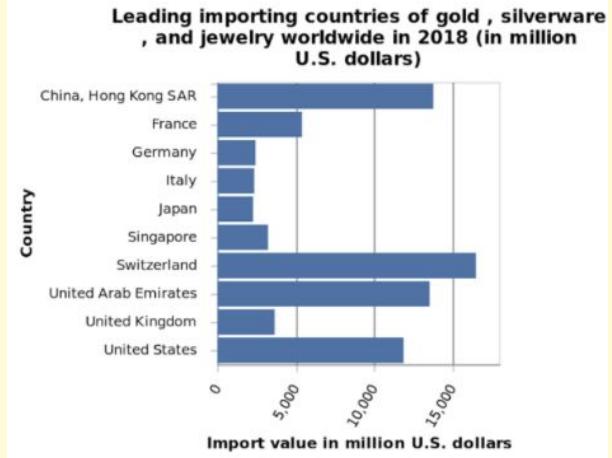
No Visual Input

Answer: No.

GPT-4V: As of my last training cut-off in January 2022, Hong Kong SAR (Special Administrative Region) has traditionally been a significant importer of gold, silverware, and jewelry, serving as a major hub for jewelry and precious metals trade. Hong Kong's role as a trade hub means that it often re-exports a considerable amount of what it imports. Specific rankings for importing countries change from year to year based on various economic factors. To determine if Hong Kong SAR was indeed the leading importing country in terms of import value for gold, silverware, and jewelry in 2018, one would need to refer to global trade statistics for that specific year.

LLaVA-1.5: According to the World Bank's Comtrade database, China, Hong Kong SAR was the leading importing country of gold, silverware, and jewelry in 2018, with an **import value of US\$157.4 billion**. This accounted for approximately 17.5% of the global total for these types of products. It's worth noting that Hong Kong SAR is a special administrative region of China, and as such, its trade data is often closely linked to that of mainland China. Therefore, it's possible that the high import value for gold, silverware, and jewelry in Hong Kong SAR in 2018 reflects trade activity taking place within China as well.

Prompt: Is China, Hongkong SAR the leading importing country of gold, silverware, and jewelry with the highest import value?



Switzerland has the highest import value.

Answer: No.

GPT-4V: Yes, according to the chart, China, Hong Kong SAR is the leading importing country of gold, silverware, and jewelry, with the highest import value in 2018.

LLaVA-1.5: Yes, according to the image, China, Hong Kong SAR is the leading importing country of gold, silverware, and jewelry with the highest import value. The chart shows that the country imported gold, silverware, and jewelry worth over 25 billion U.S. dollars in 2018. This indicates that China, Hong Kong SAR is a significant market for these products, and it plays a crucial role in the global trade of gold, silverware, and jewelry.

Figure 16. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. Sometimes, the question might **not have an affirmative answer without the context**, which GPT-4V can handle very well by **giving an uncertain answer without hallucination**. However, given the chart with all information as the input, GPT-4V and LLaVA-1.5 are unable to extract correct information to give the right answer.

Visual Supplement – Chart

Prompt: From 2001 to 2019, Is 2010 exported the most value of bread, biscuits, wafers, cakes, and pastries from the UK?

No Visual

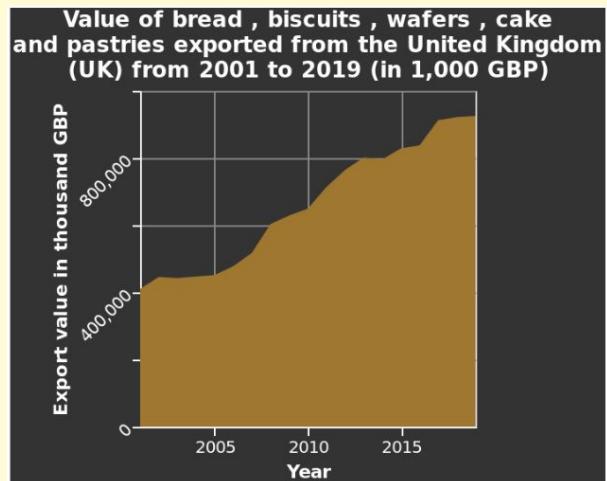
Answer: No.

GPT-4V: To answer this question accurately, one would typically refer to trade and export statistics from the UK government or relevant agencies that track exports of specific goods, including items like bread, biscuits, wafers, cakes, and pastries. However, as of my last update in January 2022, I don't have a specific year-by-year breakdown of exports for these goods from the UK from 2001 to 2019.

If you're looking for this specific information, I'd recommend checking sources such as the UK's Office for National Statistics (ONS) or similar organizations that provide detailed trade statistics.

LLaVA-1.5: According to the United Kingdom's Office for National Statistics (ONS), the UK has been a net importer of food and beverages, including bakery products, since 2001. This means that the UK has imported more of these products than it has exported. In 2019, the UK exported around £2.3 billion of bakery products while importing around £3.6 billion. Therefore, it is unlikely that any single year, such as 2010, would have been the year with the highest export value of these products from the UK.

Prompt: From 2001 to 2019, Is 2010 exported the most value of bread, biscuits, wafers, cakes and pastries from the UK?



2019/2020 has the highest export value.

Answer: No.

GPT-4V: No, according to the graph, the value of exports for bread, biscuits, wafers, cakes, and pastries from the UK increased over the years, and the highest value was reached in 2019. 2010 did not have the highest export value during this period.

LLaVA-1.5: Yes, the image shows that **in 2010, the UK exported the most value of bread, biscuits, wafers, cakes, and pastries, with a total value of 1,000 GBP.** This indicates that 2010 was a significant year for the export of these baked goods from the UK.

Figure 17. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. GPT-4V can handle uncertain queries by **giving an uncertain answer without hallucination**. Given a chart as the input, GPT-4V is able to answer the question correctly. In contrast, LLaVA-1.5 can answer the question well without the image context; however, with the addition of the image, confusion arises, and its answer becomes unclear and heavily guided by the way the question is framed due to limited vision capability. Additionally, hallucination appeared in LLaVA's answer (i.e., 1,000 GBP), which is extracted from the title of the chart, and 1,000 GBP means the total number from 2001 to 2019. We attribute this failure to the lack of visual reasoning ability.

Visual Supplement – Chart

Prompt: From 2000-2050, is the population aged between 15 and 29 years in China the most in 2000?

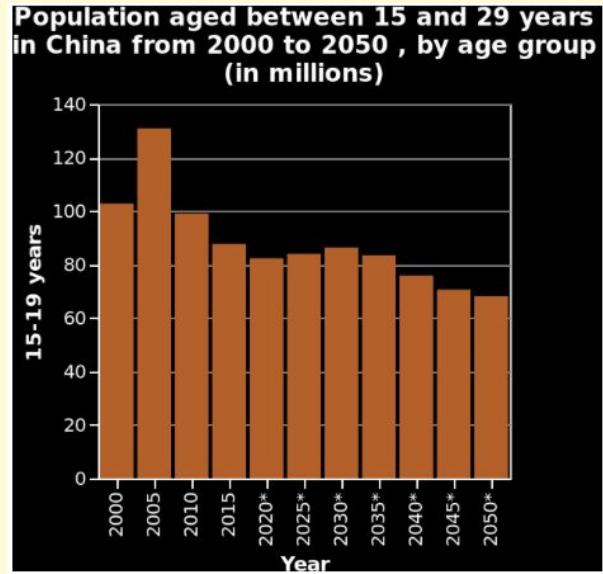
No Visual

Answer: No.

GPT-4V: As of my last update in January 2022, I don't have population projections for China up to 2050. However, I can provide some insight into the trends leading up to 2022 and offer some general expectations based on that.

LLaVA-1.5: Yes, the population aged between 15 and 29 years in China was the most in 2000. According to the United Nations Population Division, the population in this age group in China peaked in 2000 at around 306 million people. This was due to the large cohort of baby boomers born in the 1960s and 1970s, who reached their peak reproductive years in the 1980s and 1990s. However, since then, the population in this age group has been declining due to the country's aging population and declining fertility rates.

Prompt: From 2000-2050, is the population ages between 15 and 29 years in China the most in 2000?



2005 has most population from the group.

Answer: No.

GPT-4V: Yes, based on the provided graph, the population aged between 15 and 29 years in China was the highest in the year 2000.

LLaVA-1.5: Yes, the population aged between 15 and 29 years in China was the most in 2000. The image shows a bar graph with the population in this age group peaking in 2000 and then declining over the next few years. This indicates that the population in this age group was the highest in 2000 and has been decreasing since then.

Figure 18. We highlight the incorrect answer according to **visual illusion**, **language hallucination**, or **potentially mixed**. GPT-4V can handle uncertain queries by **giving an uncertain answer without hallucination**. However, given a chart as the input, GPT-4V is unable to understand it correctly. LLaVA-1.5 tends to generate answers based on the parametric memory since it generates the same answers with or without the image context.

Visual Supplement – Map

Prompt: Is New Mexico state larger than Texas state?

No Visual

Answer: No.

GPT-4V: No, New Mexico is not larger than Texas. Texas is the second-largest state in the U.S. by land area, following only Alaska. New Mexico is the fifth-largest state by land area. Texas has a land area of about 268,596 square miles, while New Mexico covers about 121,589 square miles. On any map that accurately represents the area, Texas will always appear larger than New Mexico.

LLaVA-1.5: No, New Mexico is not larger than Texas. Texas is the second-largest state in the United States in terms of both area and population, while New Mexico is the fifth-largest state in terms of area but the 36th-largest in terms of population.

Prompt: Is New Mexico state larger than Texas on the map?



New Mexico and Texas are swapped in the following image.

Answer: Yes.

GPT-4V: No, based on the map you provided, Texas is larger than New Mexico.

LLaVA-1.5: No, from the map, Texas state is larger than New Mexico state.

Prompt: Does the prime meridian intersect the Atlantic Ocean, Southern Ocean, and Arctic Ocean?

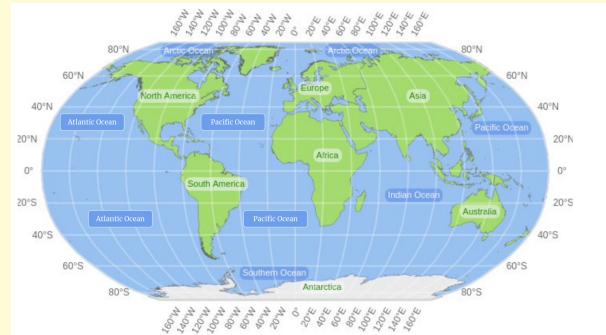
No Visual

Answer: Yes.

GPT-4V: Yes, the Prime Meridian, which is the line of 0° longitude, does intersect all three of these oceans: the **Atlantic Ocean**: The Prime Meridian runs from the North Pole southward through the Atlantic Ocean, passing west of Europe and Africa, before reaching the Southern Ocean. **Southern Ocean**: As the Prime Meridian continues south from the Atlantic, it intersects the Southern Ocean and continues all the way to the South Pole. **Arctic Ocean**: At its northernmost extent, the Prime Meridian intersects the Arctic Ocean as it goes from the North Pole southward.

LLaVA-1.5: The prime meridian, which is the line of longitude that passes through Greenwich, England, intersects the Atlantic Ocean, the Southern Ocean, and the Arctic Ocean. The prime meridian passes through the Atlantic Ocean to the west of the Azores, through the Southern Ocean south of Africa, and through the Arctic Ocean north of Canada.

Prompt: According to the map, does the prime meridian intersect the Atlantic Ocean, Southern Ocean, and Arctic Ocean?



Atlantic Ocean is replaced with the Pacific Ocean in the following image.

Answer: No, the prime meridian intersect the Pacific Ocean, Southern Ocean and arctic Ocean in this figure.

GPT-4V: Yes, according to the map, the prime meridian intersects the Atlantic Ocean, Southern Ocean, and the Arctic Ocean.

LLaVA-1.5: Yes, the prime meridian intersects the Atlantic Ocean, Southern Ocean, and Arctic Ocean.

Figure 19. We highlight the incorrect answer according to **visual illusion**, **language hallucination**, or **potentially mixed**. Both GPT-4V and LLaVA-1.5 have **an affirmative and correct answer without context**. Given an edited map (contradictory to common sense), GPT-4V and LLaVA-1.5 don't understand the map correctly. GPT-4V and LLaVA-1.5 tend to generate answers based on their parametric memory instead of the actual visual context, even though the way they answer the question indicates that the reply is based on the provided map.

Visual Supplement – Map

Prompt: Among Arizona, Oklahoma, North Carolina, and Delaware, Is Delaware the farthest north?

No Visual

Answer: Yes, Delaware is the farthest north among Arizona, Oklahoma, North Carolina, and Delaware.

GPT-4V: Yes, among the states you listed—Arizona, Oklahoma, North Carolina, and Delaware, **Delaware is the farthest north.**

LLaVA-1.5: No, Delaware is not the farthest north among Arizona, Oklahoma, North Carolina, and Delaware. **North Carolina is the farthest north among these four states,** Delaware is located to the north of Maryland and to the east of Pennsylvania, while Arizona is located to the southwest of Delaware, Oklahoma is located to the west of Delaware, and North Carolina is located to the north of Delaware.

Prompt: Among Arizona, Oklahoma, North Carolina, and Delaware on the map, Is Delaware the farthest north?



Delaware and Arizona are swapped in the following image.

Answer: No, Arizona is the farthest north among Arizona, Oklahoma, North Carolina, and Delaware.

GPT-4V: Yes, among Arizona, Oklahoma, North Carolina, and Delaware on the map, **Delaware is the farthest north.**

LLaVA-1.5: Yes, **Delaware is the farthest north among the four states on the map.**

Figure 20. We highlight the incorrect answer according to **visual illusion**, **language hallucination**, or **potentially mixed**. When there is no image context, both GPT-4V and LLaVA-1.5 have **an affirmative and correct answer**, but only GPT-4V provides the correct answer. We conclude that GPT-4V is more knowledgeable than LLaVA-1.5. With an altered map as a visual input, both models fail to reason the location relationship among the four states due to strong parametric memory.

Visual Supplement – Optical Character Recognition

Prompt: Does π range from 3.1415926 to 3.1415927?

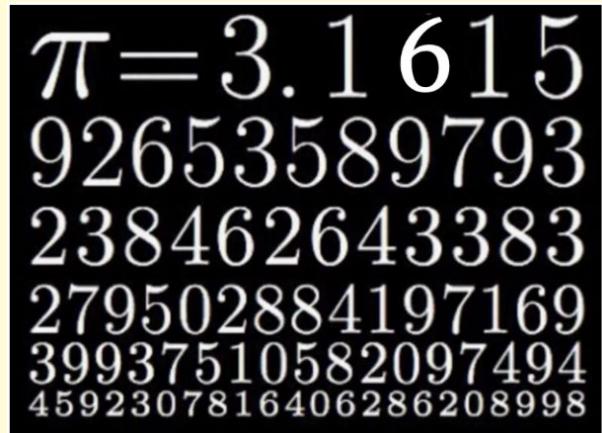
No Visual

Answer: Yes.

GPT-4V: Yes, according to the image, the value of π starts with 3.1415926, which is within the range of 3.1415926 to 3.1415927.

LLaVA-1.5: No, π (pi) is an irrational number, which means it cannot be expressed as a finite decimal or fraction. It is an infinite, non-repeating sequence of digits. The value of π is approximately 3.14159265358979323846, but it is not exactly equal to any finite decimal or fraction.

Prompt: Does π range from 3.1415926 to 3.1415927 according to the image?



$\pi = 3.141592653589793238462643383279502884197169399375105820974944592307816406286208998$

3.1415 is changed into 3.1615 in the following image.

Answer: No.

GPT-4V: Yes, according to the image, the value of π starts with 3.1415926, which is within the range of 3.1415926 to 3.1415927.

LLaVA-1.5: Yes, the image shows that π ranges from 3.1415926 to 3.1415927.

Figure 21. We highlight the incorrect answer according to visual illusion, language hallucination, or potentially mixed. When there is no image context, GPT-4V has an affirmative answer, while LLaVA-1.5 provides an irrelevant answer that is technically correct due to failed text understanding. With an altered π value as a visual input, both models fail to interpret the value from the image.

Visual Supplement – Chart

Prompt: Does China have the largest population in the world, followed by India and then the USA?

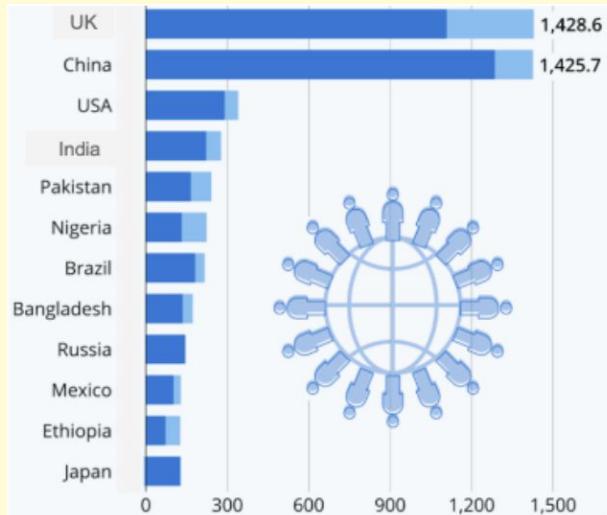
No Visual

Answer: Yes.

GPT-4V: As of my last knowledge update in September 2021, China indeed had the largest population in the world, followed by India and then the United States. However, population figures can change over time due to factors like birth rates, death rates, and immigration. It's possible that these rankings have changed since then. To get the most up-to-date population statistics, I recommend checking a reliable source like the United Nations or the World Bank, which regularly updates and publishes demographic data for countries around the world.

LLaVA-1.5: Yes, that is correct. According to the latest estimates, China has the largest population in the world, with over 1.4 billion people. India is the second most populous country, with over 1.3 billion people, and the United States is the third most populous country, with over 331 million people.

Prompt: Does China have the largest population in the world, followed by India and then the USA?



UK is the largest one, followed by China and USA.

Answer: No.

GPT-4V: Yes, according to the chart, China has the largest population, followed by India and then the USA.

LLaVA-1.5: Yes, according to the image, China has the largest population in the world, followed by India and then the USA. The graph shows the population of various countries, with China at the top, India in second place, and the USA in third. This information is based on the data available at the time the graph was created, and it may not reflect the current population statistics.

Figure 22. We highlight the incorrect answer according to visual illusion, language hallucination, or potentially mixed. Both GPT-4V and LLaVA-1.5 have an affirmative and correct answer without context, and the answer provided by GPT-4V is more rigorous since it also emphasizes the time-dependency of this knowledge. However, both models fail to answer the question correctly based on the data provided in the chart.

Visual Supplement – Table

Prompt: In the 2008 Beijing Olympics, Does China have the most gold metals, followed by the USA and Russia?

No Visual

Answer: Yes.

GPT-4V: Yes, during the 2008 Beijing Olympics, **China topped the gold medal count. They won a total of 51 gold medals.** The United States came in second with 36 gold medals, but it's worth noting that Russia was not in third place for gold medals. Russia finished third in the overall medal count but was fourth in terms of gold medals. Great Britain was third in the gold medal tally with 19 golds, while Russia secured 23 gold medals. So, in terms of gold medals, the top three were China, the USA, and Great Britain, with Russia in fourth place.

LLaVA-1.5: Yes, that is correct. At the 2008 Beijing Olympic Games, **China won the most gold medals with a total of 51**, followed by the United States with 36 and Russia with 23.

Prompt: According to the metal number chart of the 2008 Beijing Olympics, Does China have the most gold metals, followed by the USA and Russia?

	Gold	Silver	Bronze
China	23	21	28
United States	36	38	36
Russian Fed.	23	21	28
Great Britain	19	13	15
Germany	16	10	15
Australia	14	15	17
Korea	13	10	8
Japan	9	6	10
Italy	8	10	10
France	7	16	17

The gold medal number of China is changed to 23, ranked second.

Answer: No.

GPT-4V: Based on the chart provided, **China has 36 gold medals. The United States has 19 gold medals.** The Russian Federation (Russian Fed.) has 23 gold medals. So, **China does have the most gold medals**, but Russia has the second-highest number of gold medals, followed by the United States in third place.

LLaVA-1.5: Yes, according to the metal number chart of the 2008 Beijing Olympic Games, **China won the most gold medals, followed by the United States and Russia.** The chart shows the number of gold, silver, and bronze medals won by each country during the games.

Figure 23. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. When there is no visual input, both GPT-4V and LLaVA-1.5 have **an affirmative and correct answer**. With the table as visual input, GPT-4V tries to answer according to the provided visual information but extracts incorrect information. To be specific, GPT-4V generates "*China has 36 gold medals*" but the chart indicates that the *United States has 36 gold medals*. LLaVA-1.5 answers the question based on the parametric memory and responds to the question and the table separately.

Visual Supplement – Chart

Prompt: Is China the third largest country by land area in the world?

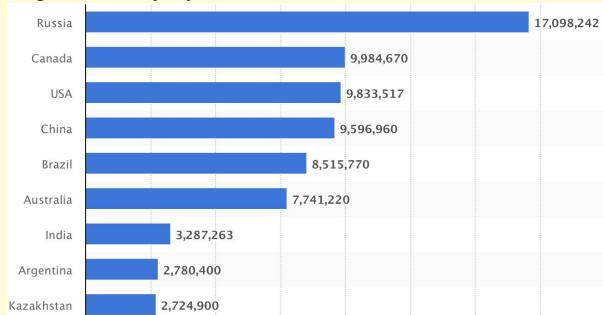
No Visual

Answer: Yes, China is the third largest country.

GPT-4V: No, China is not the third largest country by land area in the world; **it is the third largest by total area (which includes both land and water)**. By land area alone, China is the second largest.

LLaVA-1.5: Yes, **China is the third largest country in the world by land area**. It has a total land area of approximately 9.6 million square kilometers. The country is bordered by 14 other countries, including Russia, India, and Vietnam. China is also home to a diverse range of landscapes, including mountains, deserts, and forests.

Prompt: According to the chart, is China the third largest country by land area in the world?



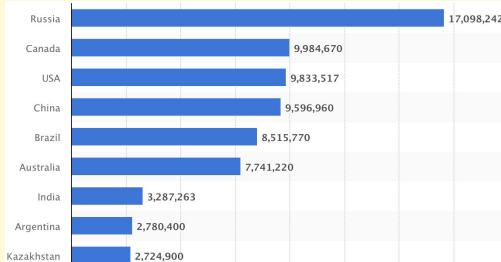
The land size of China and USA are swapped.

Answer: No, China is the fourth in the chart.

GPT-4V: No, according to the chart provided, China is the fourth largest country by land area in the world. The top three are Russia, Canada, and the USA.

LLaVA-1.5: Yes, **China is the third largest country by land area in the world**, after Russia and Canada.

Prompt: Is China the third largest country by land area in the world?



Original image

Answer: No, China is the fourth in the chart.

GPT-4V: No, according to the chart provided, China is the fourth largest country by land area in the world. The top three are Russia, Canada, and the USA.



Edited: rotated chart.

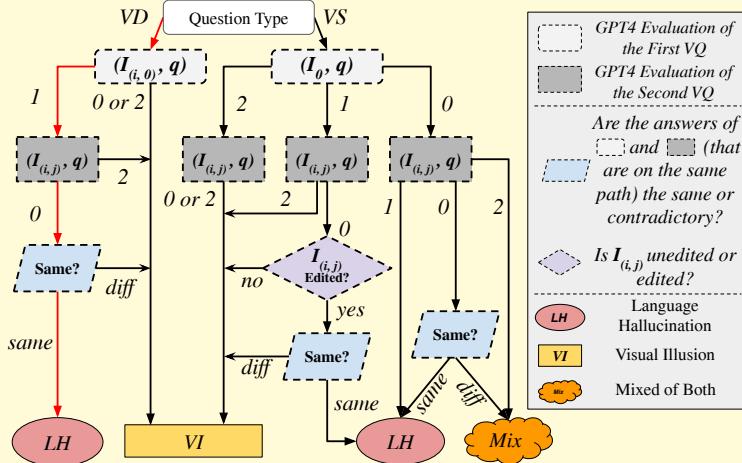
Answer: No, China is the fourth in the chart.

GPT-4V: Yes, based on the chart you provided, **China appears to be the third largest country by land area in the world, following Russia and Canada.**

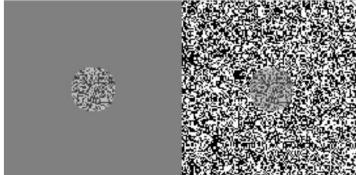
Figure 24. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. **Top:** When there is no image context, both GPT-4V and LLaVA-1.5 have **an affirmative and correct answer**, but only GPT-4V provides the correct and most rigorous answer. With the chart as visual input, GPT-4V can follow the information provided in the chart faithfully, while LLaVA-1.5 tries to answer the question based on the parametric memory. **Bottom:** We find that the predictions of GPT-4V changed completely after we flipped the chart. In contrast to the chart at the top, this failure can be attributed to the visual illusion.

Visual Dependent – Decision Tree Logic

Decision Tree: The RED path means the decision path.



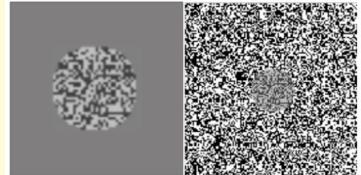
Prompt: The image shows the famous Chubb illusion. Are the center areas of two rectangular fields identical?



Original image

Answer: Yes.

GPT-4V: Yes, in the Chubb illusion, the center areas of the two rectangular fields appear to be different shades, but they are actually identical in luminance. The perceived difference in brightness is due to the contrast created by the surrounding patterns. This illusion highlights how context can influence our perception.



Edited: Circle on the left are enlarged.

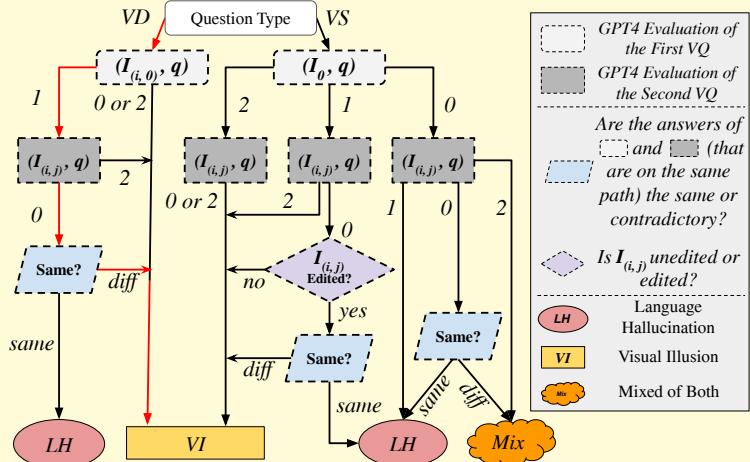
Answer: No, they are not identical.

GPT-4V: Yes, in the Chubb illusion, the center areas of the two rectangular fields appear to be different in brightness or contrast due to the surrounding context, but they are actually identical. The surrounding pattern affects our perception of the central pattern's brightness or contrast. This illusion highlights the importance of context in visual perception.

Figure 25. It is a visual-dependent sample (VD). The answer regarding the original image is correct (1), but the answer to the edited image is incorrect (0), and the two answers are the same. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Language Hallucination*. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**.

Visual Dependent – Decision Tree Logic

Decision Tree: The RED path means the decision path.



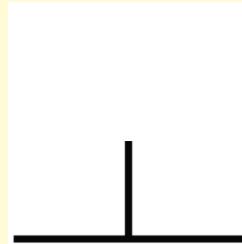
Prompt: Is the vertical line the same length as the horizontal line?



Original: The lines have the same length.

Answer: Yes, the lines have the same length.

GPT-4V: Yes, the vertical line is the same length as the horizontal line.



Edited: The vertical line is shortened.

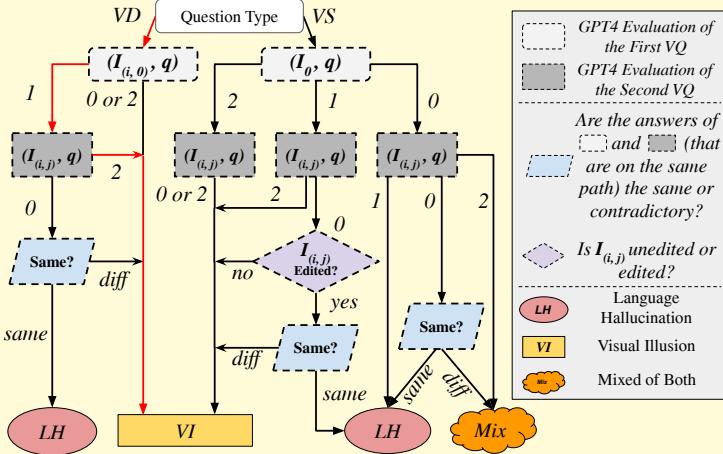
Answer: No, the vertical line is shorter.

GPT-4V: No, the vertical line appears to be longer than the horizontal line in the image provided.

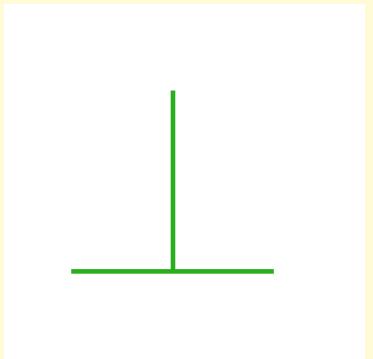
Figure 26. It is a visual-dependent sample (VD). The answer regarding the original image is correct (1), but the answer to the edited image is incorrect (0), and the two answers are not the same. This shows that GPT-4V can not compare the length of the two lines correctly. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Visual Illusion*. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**.

Visual Dependent – Decision Tree Logic

Decision Tree: The RED path means the decision path.



Prompt: Is the vertical line longer than the horizontal line?



Original image

Answer: No, the lines have the same length.

GPT-4V: The vertical line is not longer than the horizontal line.



Edited: The vertical line is elongated.

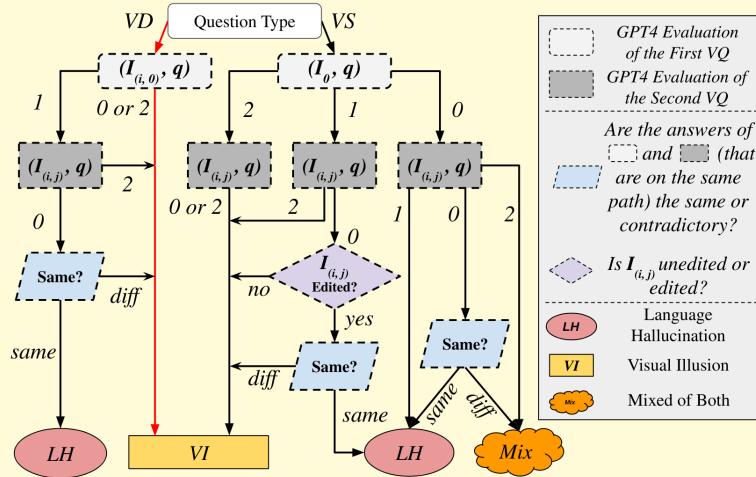
Answer: Yes, the vertical line is longer.

GPT-4V: Without the ability to measure the lines, and based on common visual illusions, the lines may be of the same length. However, this type of image is often used to demonstrate the vertical-horizontal illusion, where the vertical line appears longer than the horizontal line due to the way our perception works, even though they are the same length. If this image is not intended to be an illusion and without exact measurement, I cannot definitively say whether one line is longer than the other.

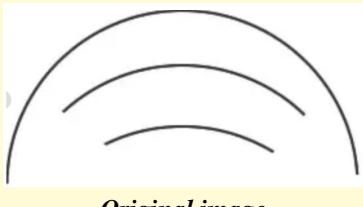
Figure 27. It is a visual-dependent sample (VD). The answer regarding the original image is correct (1), but the answer to the edited image is uncertain (2). This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Visual Illusion*. We highlight the **uncertain answer**, and the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. In this case, our decision tree can only detect *visual illusion*; we can get additional information about the failure by examining the responses.

Visual Dependent – Decision Tree Logic

Decision Tree: The RED path means the decision path.



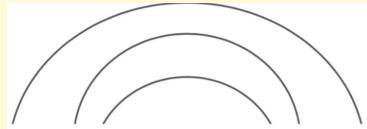
Prompt: According to the Torsion dystonia illusion, are the diameters of the three circles the same?



Original image

Answer: Yes, the curves have the same diameters.

LLaVA-1.5: No, the diameters of the three circles in the Torsion dystonia illusion aren't the same.



Edited: The curves have different diameters.

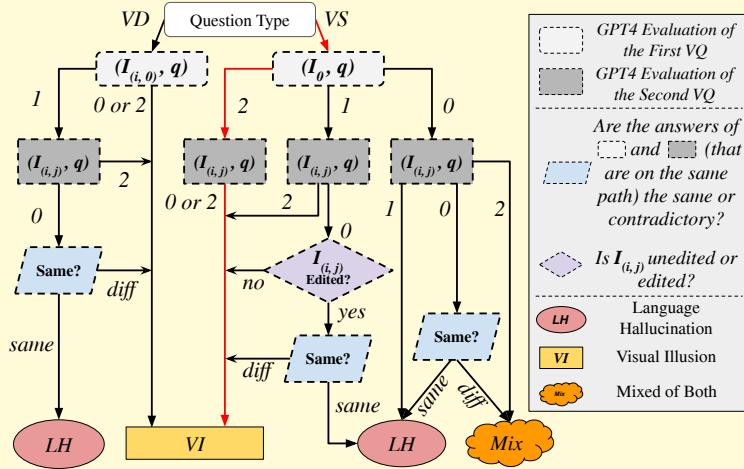
Answer: No, the curves have different diameters.

LLaVA-1.5: No, the diameters of the three circles in the Torsion dystonia illusion are not the same. The illusion is designed to make it appear as if the circles are of the same size, but they are actually different in size.

Figure 28. It is a visual-dependent sample (VD). The answer regarding the original image is incorrect (0). This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Visual Illusion*. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**.

Visual Supplement – Decision Tree Logic

Decision Tree: The RED path means the decision path.



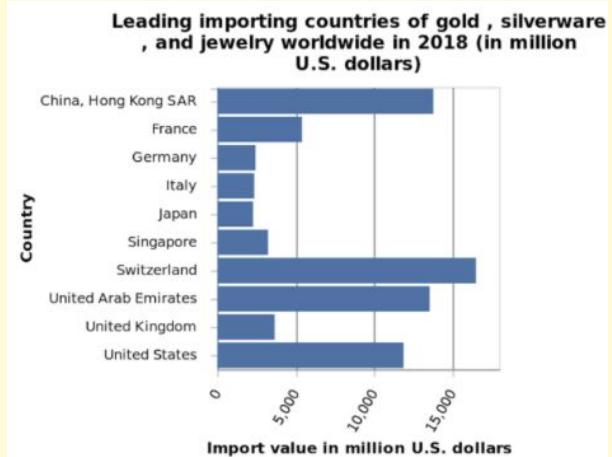
Prompt: Is China, Hongkong SAR, the leading importing country of gold, silverware, and jewelry with the highest import value in 2018?

No Visual Input

Answer: No.

GPT-4V: As of my last training cut-off in January 2022, Hong Kong SAR (Special Administrative Region) has traditionally been a significant importer of gold, silverware, and jewelry, serving as a major hub for jewelry and precious metals trade. Hong Kong's role as a trade hub means that it often re-exports a considerable amount of what it imports. Specific rankings for importing countries change from year to year based on various economic factors. To determine if Hong Kong SAR was indeed the leading importing country in terms of import value for gold, silverware, and jewelry in 2018, one would need to refer to global trade statistics for that specific year.

Prompt: Is China, Hongkong SAR the leading importing country of gold, silverware, and jewelry with the highest import value in 2018?



Switzerland has the highest import value.

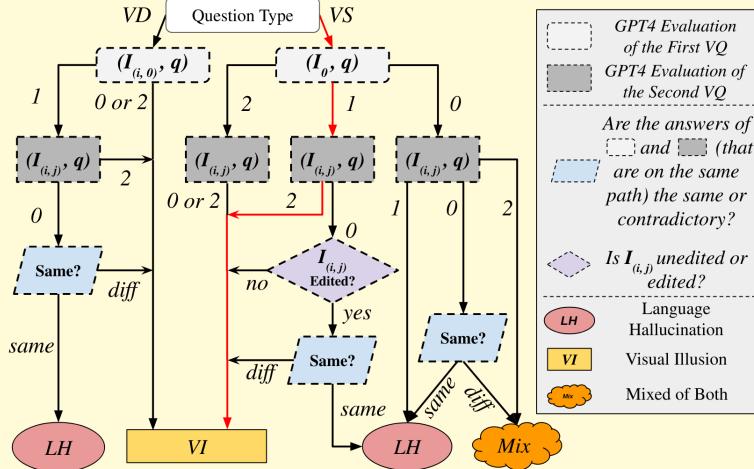
Answer: No.

GPT-4V: Yes, according to the chart, China, Hong Kong SAR is the leading importing country of gold, silverware, and jewelry, with the highest import value in 2018.

Figure 29. It is a visual-supplement sample (VS). Without an image, the answer is uncertain (2), but the answer to the edited image is incorrect (0), and the two answers are not the same. This shows that GPT-4V can not understand the chart correctly. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Visual Illusion*. We highlight the *uncertain answer*, and the *incorrect answer* according to *visual illusion*, *language hallucination*, or *potentially mixed*.

Visual Supplement – Decision Tree Logic

Decision Tree: The **RED** path means the decision path.



Prompt: Is "We shall fight on the beaches" made in the U.K.?

No Visual

Answer: Yes.

GPT-4V: The phrase "We shall fight on the beaches" is famously associated with a speech made by Winston Churchill, the Prime Minister of the United Kingdom during World War II. This speech was indeed made in the UK. Churchill delivered it to the House of Commons of the British Parliament on June 4, 1940, as part of his updates on the evolving situation in World War II, particularly the Battle of Dunkirk. The speech is renowned for its powerful rhetoric and is considered one of Churchill's most inspiring orations, designed to motivate and galvanize the British public during a time of significant adversity.

Prompt: Is the text in this figure made in the U.K.?

We shall go on to the end. We shall fight in France, we shall fight on the seas and oceans, we shall fight with growing confidence and growing strength in the air. We shall defend our Island, whatever the cost may be. We shall fight on the beaches, we shall fight on the landing grounds, we shall fight in the fields and in the streets, we shall fight in the hills. We shall never surrender, and even if, which I do not for a moment believe, this Island or a large part of it were subjugated and starving, then our Empire beyond the seas, armed and guarded by the British Fleet, would carry on the struggle, until, in God's good time, the New World, with all its power and might, steps forth to the rescue and the liberation of the old.

Original image

Answer: Yes, the text in this figure is made in the U.K.

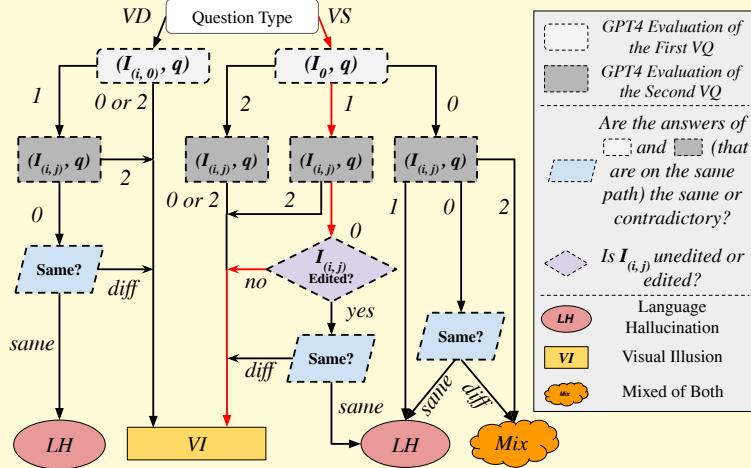
GPT-4V: The text in the image you've provided is indeed associated with the United Kingdom, as it is an excerpt from a speech by Winston Churchill, who was the Prime Minister of the United Kingdom during World War II. The speech was delivered in the House of Commons and is a part of British historical documentation.

However, whether the specific rendering of the text in the image you've provided was made in the U.K. cannot be determined from the text alone. The creation of such an image could be done anywhere in the world.

Figure 30. It is a visual-supplement sample (VS). The answer is correct (1) without being given any image. However, the answer to the edited image is uncertain (2). This corresponds to the **RED** route in the decision tree, leading to the diagnostic result of *Visual Illusion*. We highlight the **uncertain answer**, and the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**.

Visual Supplement – Decision Tree Logic

Decision Tree: The RED path means the decision path.



Prompt: Is "Ask Not What Your Country Can Do For You..." made in Canada?

No Visual

Answer: No, "Ask Not What Your Country Can Do For You..." is not made in Canada..

LLaVA-1.5: No, the quote "Ask not what your country can do for you, ask what you can do for your country" was not made in Canada. It was delivered by John F. Kennedy, the 35th President of the United States, during his inaugural address on January 20, 1961, at the United States Capitol.

Prompt: Is the text in this figure made in Canada?

Let the word go forth from this time and place, to friend and foe alike, that the torch has been passed to a new generation of Americans—born in this century, tempered by war, disciplined by a hard and bitter peace, proud of our ancient heritage—and unwilling to witness or permit the slow undoing of those human rights to which this nation has always been committed, and to which we are committed today at home and around the world.

And so, my fellow Americans: ask not what your country can do for you—ask what you can do for your country.

Original image

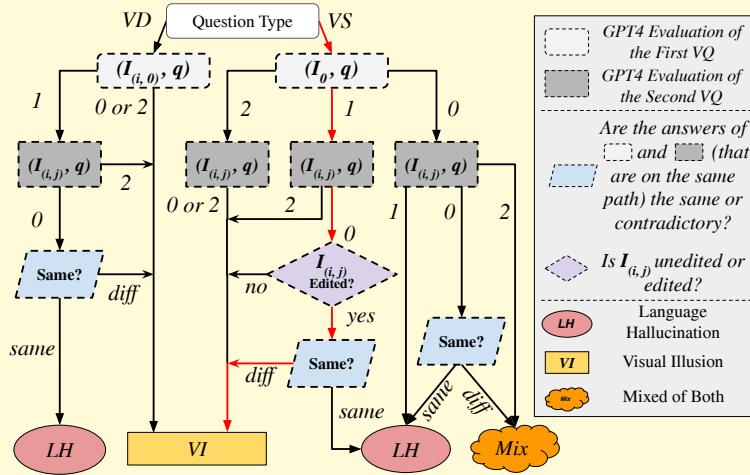
Answer: No, the text in this figure is not made in Canada.

LLaVA-1.5: Yes, **the text in this figure is made in Canada.**

Figure 31. It is a visual-supplement sample (VS). The answer is correct (1) without being given any image. However, the answer to the edited image is incorrect (0) and the two answers are not the same. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Visual Illusion*. We highlight the **incorrect answer** according to **Visual Illusion**, **language hallucination**, or **potentially mixed**.

Visual Supplement – Decision Tree Logic

Decision Tree: The RED path means the decision path.



Prompt: In the 2008 Beijing Olympics, Does China have the most gold metals, followed by the USA and Russia?

No Visual

Answer: Yes.

GPT-4V: Yes, during the 2008 Beijing Olympics, China topped the gold medal count. They won a total of 51 gold medals. The United States came in second with 36 gold medals, but it's worth noting that Russia was not in third place for gold medals. Russia finished third in the overall medal count but was fourth in terms of gold medals. Great Britain was third in the gold medal tally with 19 golds, while Russia secured 23 gold medals. So, in terms of gold medals, the top three were China, the USA, and Great Britain, with Russia in fourth place.

Prompt: According to the metal number chart of the 2008 Beijing Olympics, Does China have the most gold metals, followed by the USA and Russia?

	Gold	Silver	Bronze
China	23	21	28
United States	36	38	36
Russian Fed.	23	21	28
Great Britain	19	13	15
Germany	16	10	15
Australia	14	15	17
Korea	13	10	8
Japan	9	6	10
Italy	8	10	10
France	7	16	17

The gold medal number of China is changed to 23, ranked second.

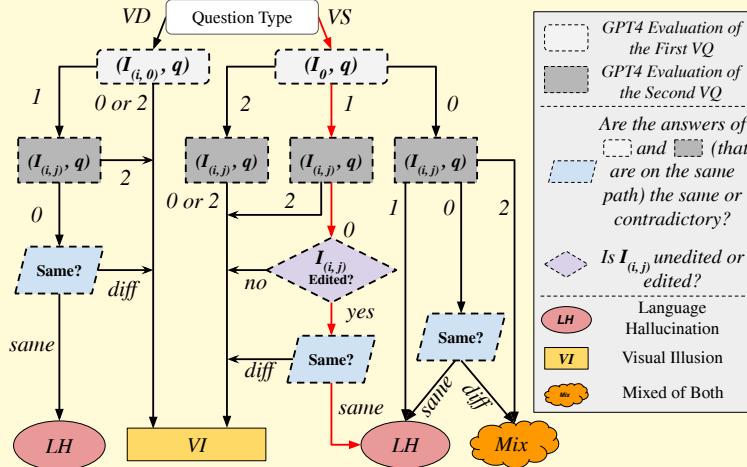
Answer: No.

GPT-4V: Based on the chart provided, **China has 36 gold medals**. The United States has 19 gold medals. The Russian Federation (Russian Fed.) has 23 gold medals. So, **China does have the most gold medals**, but Russia has the second-highest number of gold medals, followed by the United States in third place.

Figure 32. It is a visual-supplement sample (VS). The answer is correct (1) without being given any image. However, the answer to the edited image is incorrect (0). The image is edited and the two answers are not the same. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Visual Illusion*. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**. In this case, our decision tree can only detect **visual illusion**; we can get additional information about the failure by examining the responses.

Visual Supplement – Decision Tree Logic

Decision Tree: The RED path means the decision path.



Prompt: Among Arizona, Oklahoma, North Carolina, and Delaware on the map, Is Delaware the farthest north?

No Visual

Answer: Yes, Delaware is the farthest north among Arizona, Oklahoma, North Carolina, and Delaware.

GPT-4V: Yes, among the states you listed—Arizona, Oklahoma, North Carolina, and Delaware, Delaware is the farthest north.

Prompt: Among Arizona, Oklahoma, North Carolina, and Delaware on the map, Is Delaware the farthest north?



Delaware and Arizona are swapped in the following image.

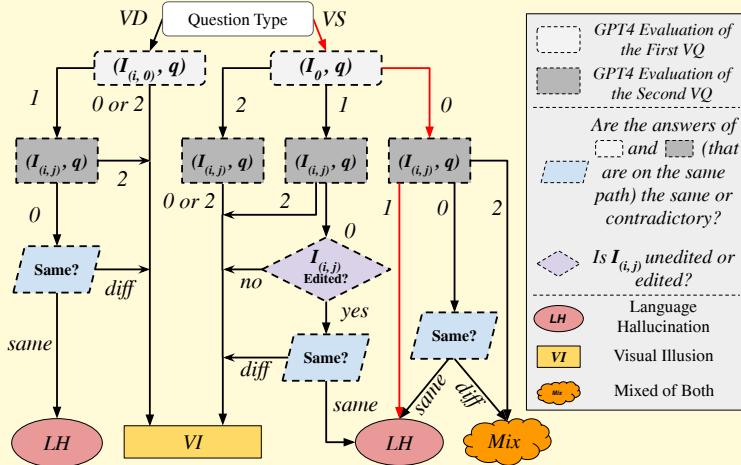
Answer: No, Arizona is the farthest north among Arizona, Oklahoma, North Carolina, and Delaware.

GPT-4V: Yes, among Arizona, Oklahoma, North Carolina, and Delaware on the map, **Delaware is the farthest north.**

Figure 33. It is a visual-supplement sample (VS). The answer without an image is correct (1), but the answer to the edited image is incorrect (0), and the two answers are the same. This shows that GPT-4V knows that "Delaware is the farthest north" in its parametric knowledge but can not answer according to the edited map. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Language Hallucination*. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**.

Visual Supplement – Decision Tree Logic

Decision Tree: The RED path means the decision path.



Prompt: Is New York the only U.S. state that shares a border with Lake Huron?

No Visual

Answer: No, New York does not share a border with Lake Huron. The only U.S. state that borders Lake Huron is Michigan. Lake Huron is one of the five Great Lakes of North America and is shared between the province of Ontario in Canada and the state of Michigan in the United States.

LLaVA-1.5: No, **New York is not the only U.S. state that shares a border with Lake Huron. In addition to New York**, the state of Michigan also shares a border with Lake Huron. Lake Huron is one of the five Great Lakes in North America, **and it borders the states of Michigan and New York on the east and west, respectively.** The other states that border Lake Huron are Ontario and Quebec in Canada.

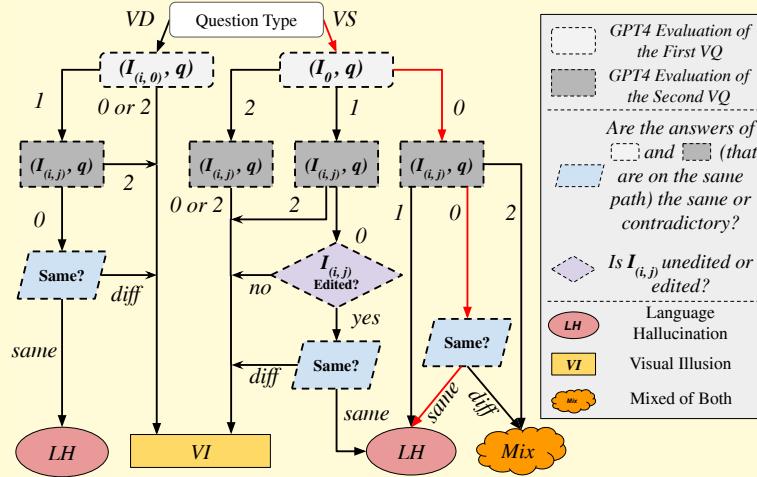
Prompt: According to the image, is New York the only U.S. state that shares a border with Lake Huron?



Figure 34. It is a visual-supplement sample (VS). The answer is incorrect (0) without being given any image. However, the answer to the original image is correct (1). This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Language Hallucination*. We highlight the **incorrect answer** according to **Language Hallucination**, **language hallucination**, or **potentially mixed**.

Visual Supplement – Decision Tree Logic

Decision Tree: The RED path means the decision path.



Prompt: From 2000-2050, is the population aged between 15 and 29 years in China the most in 2000?

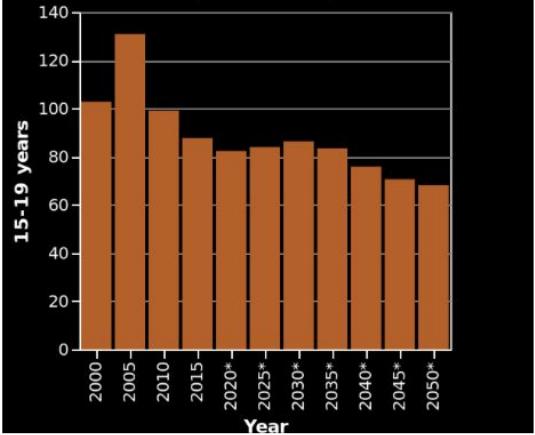
No Visual

Answer: No.

LLaVA-1.5: Yes, **the population aged between 15 and 29 years in China was the most in 2000.** According to the United Nations Population Division, the population in this age group in China peaked in 2000 at around 306 million people. This was due to the large cohort of baby boomers born in the 1960s and 1970s, who reached their peak reproductive years in the 1980s and 1990s. However, since then, the population in this age group has been declining due to the country's aging population and declining fertility rates.

Prompt: From 2000-2050, is the population aged between 15 and 29 years in China the most in 2000?

Population aged between 15 and 29 years in China from 2000 to 2050 , by age group (in millions)



2005 has most population from the group.

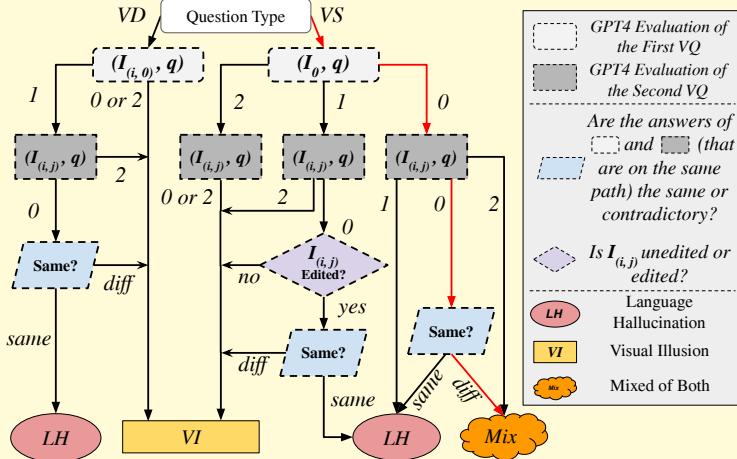
Answer: No.

LLaVA-1.5: Yes, **the population aged between 15 and 29 years in China was the most in 2000.** The image shows a bar graph with the population in this age group peaking in 2000 and then declining over the next few years. This indicates that the population in this age group was the highest in 2000 and has been decreasing since then.

Figure 35. It is a visual-supplement sample (VS). The answer is incorrect (0) without being given any image. However, the answer with a provided image is incorrect (0). And the two answers are the same. This shows that LLaVA-1.5 can not understand the rotated chart correctly. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Language Hallucination*. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**.

Visual Supplement – Decision Tree Logic

Decision Tree: The RED path means the decision path.



Prompt: Among Arizona, Oklahoma, North Carolina, and Delaware on the map, Is Delaware the farthest north?

No Visual

Answer: Yes, Delaware is the farthest north among Arizona, Oklahoma, North Carolina, and Delaware.

LLaVA-1.5: No, Delaware is not the farthest north among Arizona, Oklahoma, North Carolina, and Delaware. **North Carolina is the farthest north among these four states,** Delaware is located to the north of Maryland and to the east of Pennsylvania, while Arizona is located to the southwest of Delaware, Oklahoma is located to the west of Delaware, and North Carolina is located to the north of Delaware.

Prompt: Among Arizona, Oklahoma, North Carolina, and Delaware on the map, Is Delaware the farthest north?



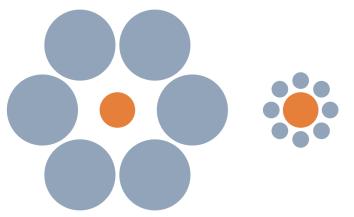
Delaware and Arizona are swapped in the following image.

Answer: No, Arizona is the farthest north among Arizona, Oklahoma, North Carolina, and Delaware.

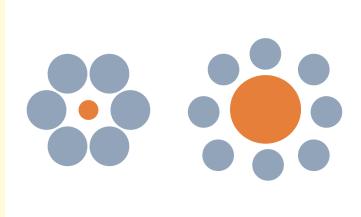
LLaVA-1.5: Yes, **Delaware is the farthest north among the four states on the map.**

Figure 36. It is a visual-supplement sample (VS). The answer is incorrect (0) without being given any image. The answer to the edited image is incorrect (0). But the two incorrect answers are not the same. This shows that the commonsense knowledge knowledge of LLaVA-1.5 is weak and wrong without the image input. Additionally, LLaVA-1.5 can not understand the map correctly through the visual input. This corresponds to the RED route in the decision tree, leading to the diagnostic result of *Potentially Mixed*. We highlight the **incorrect answer** according to **visual illusion**, **language hallucination**, or **potentially mixed**.

Gemini Pro Vision Responses



Original: Orange balls have the same size.



Edited: Orange balls have different size.

Prompt 1: Is the right orange circle the same size as the left orange circle?

Answer: Yes.

Gemini Pro: Yes.

Answer: No, right orange circle is bigger.

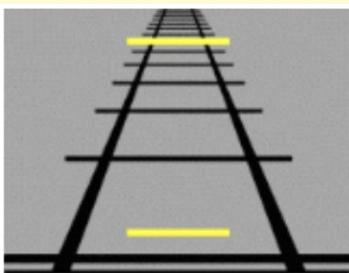
Gemini Pro: Yes.

Prompt 2: Is the right orange circle smaller than the left orange circle?

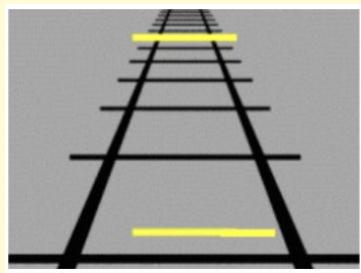
Answer: No, they are of the same size.

Gemini Pro: Yes. The right orange circle is the sun.

The left orange circle is closer to the camera.



Original: The yellow lines have the same length.



Edited: The bottom yellow line is elongated.

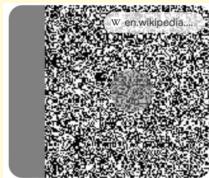
Prompt: Are the two yellow lines the same length?

Answer: Yes, the yellow lines have the same length.

Gemini Pro: Yes. The two yellow lines are the same length.

Answer: No, the bottom yellow line is longer.

Gemini Pro: Yes. The two yellow lines are the same length. They are both the same distance from the woman's hand.



Original: The line in the middle is straight.



Edited: The line in the middle is curved.

Prompt: Is the vertical line in the middle actually curved?

Answer: No, the line in the middle is straight.

Gemini Pro: No. The vertical line in the middle is not curved. It is straight.

Answer: Yes, the line in the middle is curved.

Gemini Pro: Yes, the line is curved. The line is the edge of a wave. The waves in the ocean are curved.

Figure 37. We highlight the incorrect answer.

References

- [1] Gpt-4v(ision) system card. 2023. 6, 7
- [2] Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. nocaps: novel object captioning at scale. *International Conference on Computer Vision*, pages 8947–8956, 2019. 3
- [3] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. 2, 6, 7
- [4] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015. 1
- [5] Anas Awadalla, Irena Gao, Joshua Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo, 2023. 1
- [6] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *ArXiv*, abs/2308.12966, 2023. 6, 7
- [7] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023. 1
- [8] Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Guanaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpagasus: Training a better alpaca with fewer data. *arXiv preprint arXiv:2307.08701*, 2023. 1
- [9] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, 2023. 1
- [10] Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and Huaxiu Yao. Holistic analysis of hallucination in gpt-4v(ision): Bias and interference challenges. *ArXiv*, abs/2311.03287, 2023. 4
- [11] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructclip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2023. 2, 6, 7
- [12] Danny Driess, F. Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Azyaan Wahid, Jonathan Tompson, Quan Ho Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Peter R. Flo-
- rence. Palm-e: An embodied multimodal language model. In *International Conference on Machine Learning*, 2023. 2
- [13] Zhen fei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Lu Sheng, Lei Bai, Xiaoshui Huang, Zhiyong Wang, Wanli Ouyang, and Jing Shao. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. *ArXiv*, abs/2306.06687, 2023. 3
- [14] 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. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 3, 4
- [15] Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. Multimodal-gpt: A vision and language model for dialogue with humans, 2023. 5
- [16] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. *International Journal of Computer Vision*, 127:398 – 414, 2016. 3
- [17] Tianrui Guan, Yurou Yang, Harry Cheng, Muyuan Lin, Richard Kim, Rajasimman Madhivanan, Arnie Sen, and Dinesh Manocha. Loc-zson: Language-driven object-centric zero-shot object retrieval and navigation, 2023. 1
- [18] Anish Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. *ArXiv*, abs/2308.06394, 2023. 4
- [19] MD Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. A comprehensive survey of deep learning for image captioning. *ACM Computing Surveys (CSUR)*, 51(6):1–36, 2019. 1
- [20] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *ArXiv*, abs/2307.16125, 2023. 3
- [21] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ArXiv*, abs/2301.12597, 2023. 1, 6, 7
- [22] Sheng Li and Nima Tajbakhsh. Scigraphqa: A large-scale synthetic multi-turn question-answering dataset for scientific graphs. *ArXiv*, abs/2308.03349, 2023. 4
- [23] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji rong Wen. Evaluating object hallucination in large vision-language models. *ArXiv*, abs/2305.10355, 2023. 1, 3, 4, 5
- [24] Yanda Li, Chi Zhang, Gang Yu, Zhibin Wang, Bin Fu, Guosheng Lin, Chunhua Shen, Ling Chen, and Yunchao Wei. Stablellava: Enhanced visual instruction tuning with synthesized image-dialogue data. *ArXiv*, abs/2308.10253, 2023. 2
- [25] Zongxia Li, Paiheng Xu, Fuxiao Liu, and Hyemi Song. Towards understanding in-context learning with contrastive demonstrations and saliency maps. *arXiv preprint arXiv:2307.05052*, 2023. 1

- [26] Chen Liang, Jiahui Yu, Ming-Hsuan Yang, Matthew Brown, Yin Cui, Tuo Zhao, Boqing Gong, and Tianyi Zhou. Module-wise adaptive distillation for multimodality foundation models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 1
- [27] Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. Visual news: Benchmark and challenges in news image captioning. *arXiv preprint arXiv:2010.03743*, 2020. 1, 3
- [28] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 1, 2, 3, 4, 6, 7
- [29] Fuxiao Liu, Hao Tan, and Chris Tensmeyer. Documentclip: Linking figures and main body text in reflowed documents. *arXiv preprint arXiv:2306.06306*, 2023. 1
- [30] Fuxiao Liu, Yaser Yacoob, and Abhinav Shrivastava. Covid-vts: Fact extraction and verification on short video platforms. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 178–188, 2023. 1
- [31] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 1, 5, 6, 7
- [32] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 1, 2
- [33] Yuanzhan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahu Lin. Mmbench: Is your multi-modal model an all-around player? *ArXiv*, abs/2307.06281, 2023. 3
- [34] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chun yue Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. *ArXiv*, abs/2310.02255, 2023. 4
- [35] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022. 1
- [36] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023. 2
- [37] Tanik Saikh, Tirthankar Ghosal, Amish Mittal, Asif Ekbal, and Pushpak Bhattacharyya. Scienceqa: A novel resource for question answering on scholarly articles. *International Journal on Digital Libraries*, 23(3):289–301, 2022. 1
- [38] Anthropic Team. Claude 3, 2024. 6, 7
- [39] Gemini Team. Gemini: A family of highly capable multimodal models, 2023. 6, 7
- [40] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 1
- [41] Alasdair Tran, Alexander Mathews, and Lexing Xie. Transform and tell: Entity-aware news image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13035–13045, 2020. 1
- [42] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: Lessons learned from the 2015 mscoco image captioning challenge. *IEEE transactions on pattern analysis and machine intelligence*, 39(4):652–663, 2016. 1
- [43] Bin Wang, Fan Wu, Xiao Han, Jiahui Peng, Huaping Zhong, Pan Zhang, Xiao wen Dong, Weijia Li, Wei Li, Jiaqi Wang, and Conghui He. Vigc: Visual instruction generation and correction. *ArXiv*, abs/2308.12714, 2023. 3
- [44] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. *ArXiv*, abs/2205.14100, 2022. 6, 7
- [45] Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, et al. Larger language models do in-context learning differently. *arXiv preprint arXiv:2303.03846*, 2023. 1
- [46] Yijia Xiao, Yiqiao Jin, Yushi Bai, Yue Wu, Xianjun Yang, Xiao Luo, Wenchao Yu, Xujiang Zhao, Yanchi Liu, Haifeng Chen, Wei Wang, and Wei Cheng. Large language models can be good privacy protection learners. 2023. 1
- [47] Yijun Yang, Tianyi Zhou, Kanxue Li, Dapeng Tao, Lusong Li, Li Shen, Xiaodong He, Jing Jiang, and Yuhui Shi. Embodied multi-modal agent trained by an llm from a parallel textworld, 2023. 1
- [48] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of lmms: Preliminary explorations with gpt-4v(ision), 2023. 1, 2, 5, 9
- [49] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023. 1
- [50] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023. 1, 2, 5, 6, 7
- [51] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration, 2023. 6, 7
- [52] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023. 2
- [53] Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xinguo Sun, and Enhong Chen. Woodpecker: Hallucination correction for multimodal large language models. *ArXiv*, abs/2310.16045, 2023. 3
- [54] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *ArXiv*, abs/2308.02490, 2023. 3

- [55] Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, et al. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022. 1
- [56] Yan Zeng, Hanbo Zhang, Jiani Zheng, Jiangnan Xia, Guoqiang Wei, Yang Wei, Yuchen Zhang, and Tao Kong. What matters in training a gpt4-style language model with multimodal inputs? *arXiv preprint arXiv:2307.02469*, 2023. 4
- [57] Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. Investigating the catastrophic forgetting in multimodal large language models. *arXiv preprint arXiv:2309.10313*, 2023. 1
- [58] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren’s song in the ai ocean: A survey on hallucination in large language models. *ArXiv*, abs/2309.01219, 2023. 1, 3
- [59] Yichi Zhang, Jiayi Pan, Yuchen Zhou, Rui Pan, and Joyce Chai. Grounding visual illusions in language: Do vision-language models perceive illusions like humans? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2023. 3
- [60] Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tongfei Sun. Llavar: Enhanced visual instruction tuning for text-rich image understanding. *ArXiv*, abs/2306.17107, 2023. 2
- [61] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023. 1
- [62] Kaizhi Zheng, Xuehai He, and Xin Eric Wang. Minigpt-5: Interleaved vision-and-language generation via generative vokens. *ArXiv*, abs/2310.02239, 2023. 6, 7
- [63] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 1, 2, 6, 7