Vision Language Models are Biased: Counting legs of an animal is surprisingly hard

An Vo^{1*} Khai-Nguyen Nguyen^{2*} Mohammad Reza Taesiri³ Vy Tuong Dang¹ Anh Totti Nguyen⁴ Daeyoung Kim¹

Abstract

Large language models (LLMs) memorize a vast amount of prior knowledge from the Internet that helps them on downstream tasks, but also may notoriously sway their outputs towards wrong or biased answers (Sheng et al., 2019b; Gallegos et al., 2024). In this work, we test how the knowledge about popular subjects hurts the accuracy of vision language models (VLMs) on standard, obiective visual tasks of **counting**, a common mathematical skill in everyday life. We find that stateof-the-art VLMs are strongly biased (e.g., unable to recognize a fourth stripe has been added to a 3stripe Adidas logo) scoring an average of 17.05% accuracy in counting (e.g., counting stripes in an Adidas-like logo) across 7 diverse domains from animals, logos, chess, boardgames, optical illusions, to patterned grids. Inserting text (e.g., "Adidas") describing the subject name into the counterfactual image further decreases VLM accuracy. The biases in VLMs are so strong that instructing them to double-check their results or rely exclusively on image details to answer improves counting accuracy by only +2 points, on average. Our findings reveal critical limitations of VLM capabilities in visual counting, posing an important question for how to perform math under strong perceptual bias.

1. Introduction

Large language models (LLMs) are trained on the Internet data and therefore learn a vast amount of prior knowledge that (a) help them on downstream tasks but (b) sometimes sway their answers towards wrong or biased choices (Anonymous, 2025; Sheng et al., 2019b). Interestingly, LLMs also memorize visual knowledge from its colossal text-only corpus (Sharma et al., 2024), e.g., the US national flag has 50 stars and 13 stripes or dogs have four legs (Fig. 1). Because vision language models (VLMs) are built by pre-training LLMs either exclusively on text data (i.e., for late fusion with vision encoders) (Liu et al., 2023; Bai et al., 2023) or on a mix of text, image, and multimodal data in an early fusion manner (Team, 2024), they may inherit strong biases from the text corpus when answering visual questions (Lee et al., 2023: Liu et al., 2024: Lee et al., 2025: Guan et al., 2024a).

Prior evidence (Guan et al., 2024b; Lee et al., 2025; Liu et al., 2024) showing VLMs are biased were exclusively on artificial Y/N questions that often directly contain the biased statement, e.g., "Is the mouse smaller than the cat?" (Liu et al., 2024), which is framed to contradict their counterfactual (CF) image where the cat is smaller. Therefore, it is unclear (1) how much the image contributes to VLMs' wrong answers or it is solely the textual prompt; (2) how such biases impact standard, objective visual tasks with neutral, unbiased prompts. In this work, we aim to evaluate how the knowledge of LLMs about popular subjects (e.g., dogs and the US flag) negatively impact the accuracy of VLMs on objective visual questions of object counting, identification (Q1 & Q3 in Fig. 2) and low-level visual tasks (e.g., measuring whether two lines are parallel; Fig. 1f). For example, we provide a CF image of a 3-legged chicken and ask VLMs "How many legs does this animal have?" (Fig. 1a).

Leveraging state-of-the-art (SOTA) image editors, VLMs, and image processing libraries, we propose VLMBias, a framework for automating the enumeration and generation of biased subjects, questions, and counterfactual images. We manually review all generated images and reject those that are deemed low-quality or debatable. We test VLMs on

^{*}Equal contribution ¹KAIST ²College of William and Mary ³University of Alberta ⁴Auburn University. Correspondence to: An Vo <an.vo@kaist.ac.kr>, Khai-Nguyen Nguyen <knguyen07@wm.edu>, Mohammad Reza Taesiri <mtaesiri@gmail.com>, Vy Tuong Dang <vydang@kaist.ac.kr>, Anh Totti Nguyen <anh.ng8@gmail.com>, Daeyoung Kim <kimd@kaist.ac.kr>.

The 2^{nd} AI for MATH Workshop at the 42^{nd} International Conference on Machine Learning, Vancouver, Canada. Copyright 2025 by the author(s).

Examples of VLM failures across 7 domains of VLMBias

 \mathbf{H} How many legs does this animal have? Answer with a number in curly brackets, e.g., $\{9\}$.

🐨 How many **points** are there on the star in the logo of this car? Answer with a number in curly brackets, e.g., {9}.

How many stripes are there in this flag? Answer with a number in curly brackets, e.g., {9}.

- Answer with a number in curly brackets, e.g., {9}.
- How many **rows** are there on this board? Answer with a number in curly brackets, e.g., {9}.
- Are the two horizontal lines **parallel**? Answer in curly brackets, e.g., {Yes} or {No}.
- **W** How many **circles** are there in cell C3? Answer with a number in curly brackets, e.g., {9}.

	a.	₹ €		b. 🐨		. 🗲			5 3 6 9 8 8 4 7 6 6 1	1 2 6 6 3 1 2 6 6 4 1 9 5 7 2 8 7	f.			y. III · · · · · · · · · · · · · · · · · ·
+	2	X	3	X	13	×	32	×	9	×	Yes	×	3	×
*	2	X	3	×	13	×	32	×	9	×	No	1	2	1
\$	4	X	3	×	13	×	32	X	9	×	Yes	×	3	×
(2	X	3	×	13	X	32	X	9	×	Yes	×	3	×
6	2	×	3	×	13	×	32	×	9	×	Yes	×	3	×
Bias	2	X	3	×	13	×	32	×	9	×	Yes	X	3	×
GT	3	1	4	 Image: A second s	14	1	31	√	10	 Image: A second s	No	 Image: A second s	2	1
	4	Gei	mini	-2.5 Pi	0	Sonne	t-3.7	S GI	PT-4.	1 💿	03 🚳	/ 04-	mini	

Figure 1: VLMs fail on 6 counting tasks (a-e & g) and one low-level vision task (f).

questions spanning **seven** diverse subjects in the decreasing order of popularity: (a) animals , (b) logos , (c) flags , (d) chess pieces 2; (e) game boards , (f) optical illusion , and (g) patterned grids , (see Sec. 3). For all subjects, the tasks are counting and object identification, except for optical illusion , questions, which were originally designed to test human low-level vision (e.g., identifying whether two circles are of the same size).

We test five SOTA VLMs: 3 thinking models of ★ Gemini-2.5 Pro (Google, 2025), ◎ 03 (OpenAI, 2025b) ◎ ★ 04-mini (OpenAI, 2025b); and 2 nonthinking models of ⑧ Sonnet-3.7 (Anthropic, 2025), ◎ GPT-4.1 (OpenAI, 2025a). Our key findings are:

- 1. All five VLMs recognize the VLMBias subjects from the original, unmodified image, scoring 100% accuracy on both identification and counting questions (Fig. 2a) (Sec. 4.1).
- 2. VLMs struggle to count → animal legs when one extra leg is added to 2-legged (birds; Fig. 1a) and 4-legged animals (1.01% and 2.50% accuracy, respectively; Sec. 4.2).
- 3. When logos (**) of famous car and sportswear brands are modified to have one more or one fewer of the famous visual elements (e.g., stripes on the Adidas logo; Fig. 2b), VLMs struggle to count these elements.

Their answers are extremely biased (0.44% accuracy) on CF car logos and slightly less biased on shoe logos (17.57% accuracy) (Appendix L.1). Similarly, VLMs fail to (a) detect the number of stripes and stars in the CF versions of popular flags ♥ (Appendix L.2); (b) count the chess pieces 2 chess on a chessboard when a piece is replaced or removed (Sec. 4.3); and (c) count the rows and columns of a modified board ∰ of famous board games (Appendix L.4).

- 4. On optical illusions ← (e.g., Ebbinghaus; Fig. 5), all VLMs accurately predict the names of the original, well-known illusions *but unable* to detect that the illusion graphics together with the groundtruth answers have changed, scoring around random chance (Sec. 4.4).
- 5. Unlike the above cases, we test VLMs on the patterned grid **iii** task where *no prior* information (e.g., famous illusions or logos) of the image exists on the Internet. In this grid, all cells follow a pattern except for one cell about which we will question VLMs. VLMs perform poorly, failing to detect the subtly-changed cell and answer based on the rules implied from the surrounding cells (Sec. 4.5).
- 6. To confirm VLM failures to counting (Q1 & Q2) are due to their visual bias, we further test VLMs on Y/N

identification questions (Fig. 2; Q3) but they also similarly struggle to answer (Sec. 4.6). In another experiment where the subject name (e.g., "Adidas") is added to each CF image (e.g., 4-striped logo), VLM counting accuracy further drop by -2 to -6 points, confirming the bias learned from the text corpus influences its counting (Sec. 4.7).

 Instructing VLMs to rely only on the visual details in the image alone to answer or to double-check the result in a 2nd-turn message improves their counting accuracy by at most +2 points only, confirming the severe bias of the SOTA VLMs (Sec. 4.8).

2. Related work

Bias in LLMs and VLMs LLMs exhibited biases across various domains, including social (Shin et al., 2024; Hu et al., 2025), cultural (Li et al., 2024; Naous et al., 2024; Abid et al., 2021; Wang et al., 2024), demographic (Zhao et al., 2023; Kumar et al., 2024), political (Bang et al., 2024; Potter et al., 2024), cognitive (Echterhoff et al., 2024; Koo et al., 2024), and biases related to specific names, numbers, or values (Zhang et al., 2024b; Koevering & Kleinberg, 2024). These biases often correlate with the overrepresented associations between textual cues and specific classes or attributes (e.g., associating older people with forgetfulness) (Parrish et al., 2022) in the pretraining data. Biases are not limited to textual data but extend into the visual domain. VLMs also exhibit gender biases (Hall et al., 2023; Xiao et al., 2024; Hirota et al., 2022; Fraser & Kiritchenko, 2024), stereotypical portrayals (Ruggeri & Nozza, 2023; Janghorbani & De Melo, 2023; Raj et al., 2024), and social biases (Howard et al., 2024; Sathe et al., 2024).

In our work, we investigate VLM bias in visual question answering (VQA), specifically, in cases where the visual cues in an counterfactual image strongly bias VLMs answers towards knowledge commonly known on the Internet, effectively ignoring the counterfactual (CF) modifications, resulting in inaccurate predictions (Fig. 2).

Table 1: Our VLMBias presents natural, objective counting and identification questions while prior benchmarks insert biased statements into the prompt.

Benchmark	Biased prompt	Biased image	CF images	Generation method	Adversarial injection	Top leaderboard	Question types
PhD-ccs (Liu et al., 2024)	1	×	750	DALL-E	In-prompt	GPT-40 81.2%	Y/N
VLind-Bench (Lee et al., 2025)	1	×	2,576	DALL-E	n/a	GPT-40 89.4%	Y/N
HallusionBench (Guan et al., 2024a)	1	1	181	manual	n/a	GPT-4V 31.4%	Y/N
VLMBias (ours)	×	1	1,392	automated	In-image Title	04-mini 20.25%	Counting (Q1, Q2 Y/N (Q3)

Visual Hallucination Benchmarks Many prior attempts tested VLMs on visually ambiguous images (Liu et al., 2024; Huang et al., 2024; Tong et al., 2024), optical illusion (Guan

et al., 2024a; Wu et al., 2024), CF images (Lee et al., 2025; Guan et al., 2024a) and counter-commonsense images (Liu et al., 2024; Lee et al., 2025; Bitton-Guetta et al., 2023; Zhou et al., 2023). The most relevant benchmarks (Liu et al., 2024; Lee et al., 2025; Guan et al., 2024a) (see Tab. 1) have three main drawbacks: (1) they primarily incorporate biased textual cues (e.g., directly mentioning entity names) in the questions to provoke LLM hallucination; (2) they use only Yes/No questions, which limit them to artificial questions instead of objective downstream tasks such as counting in our work; and (3) they do not study the effects of in-image adversarial injection. Among these, the Visual Dependent subset of HallusionBench (Guan et al., 2024a) is the most similar to VLMBias. However, they still have the above limitations and rely entirely on humans to manually edit images to produce 181 CF images. In contrast, we automate the CF-image generation process and humans only review the generated images.

VLMBias addresses these limitations by (1) inserting bias cues into the image keeping the prompt neutral; (2) using counting questions, which are objective and challenging to VLMs (Rahmanzadehgervi et al., 2024); and (3) optionally, injecting extra bias cues as text into the image. Furthermore, VLMBias is fully automated on 6 out of 7 tasks: LLMs to generate ideas; Python scripts to generate abstract images; and SOTA text-to-image models (Gemini-2.0 Flash & GPT-40) to produce photo-realistic images.

3. The VLMBias Benchmark

VLMBias evaluates VLMs' visual bias by presenting a pair of counting question and subtly modified versions of well-known objects (e.g., changing the Adidas logo from 3-striped to 4-striped; Fig. 2c). We choose the counting task as it is a generic, objective visual question that does not contain specific biased statements or subjects. We test whether the visual bias cues in the background is so strong that it will make VLMs default to biased answers and ignore the modifications.

Taxonomy To test VLM biases, we choose 7 unique topics of **decreasing popularity**, i.e., from common animals \succeq , then logos \boxdot to optical illusions, and a novel visual pattern **iii** that we create and did not exist before. For each topic, we generate images from scratch. The generated images are photo-realistic for 2 topics and abstract for the rest of 5 topics.

(1) Photo-realistic images are used in 2 tasks: Animals and Bogos. These images cover the most common subjects, natural (A) and man-made (B). They are created and modified by SOTA text-to-image generators (**Gemini-2.0 Flash image generation and GPT-40).

(2) Abstract images are used in 5 tasks: 🟴 flags, 🙎 chess



Figure 2: Given a subject (e.g., Adidas logo), we first confirm that *all* VLMs have sufficient knowledge about the subject via an ID and counting sanity-check questions (a). Then, we test VLMs on the counterfactual image (b) and report its accuracy on the counting (Q1 & Q2) and an Y/N identification task (Q3). For all tasks, we test the hypothesis that the visual bias cues in the **background** (c) may be so strong that it cause VLMs to ignore the modified object and default to biased answers.

Controls To minimize the language *bias* in the prompt, we use two different prompts per test image, written in neutral, descriptive terms (e.g. *stylized curves* for *Nike swooshes*). Each test image is re-scaled to $D \times D$ pixels where $D \in \{384, 768, 1152\}$. In each task, we ask 3 questions, e.g., our three questions (two counting & one identification) in the animal-leg \forall task are (Fig. 2b):

Q1: How many legs does this animal have? Answer with a number in curly brackets, e.g., {9}.

Q2: Count the legs of this animal. Answer with a number in curly brackets, e.g., {9}.

Q3: Is this an animal with 4 legs? Answer in curly brackets, e.g., {Yes} or {No}.

3.1. Task 1: Counting animal legs when an extra leg is added 🍾

Pretrained on the Internet data, VLMs must have colossal prior knowledge of the count of rainimal legs from both textual and image data. Following this hypothesis, we generate images of well-known animals but with *one extra leg* (e.g., 3-legged birds or 5-legged dogs) and ask VLMs to count legs.

Images We use a three-step image generation process: (1) We obtain a list of 100 well-known \mathcal{F} animals with two or four legs using $\mathfrak{G} \neq 04$ -mini; (2) For each \mathcal{F} animal, we ask \mathcal{F} Gemini-2.0 Flash to generate their side-view

images; (3) Then, we instruct ***** Gemini-2.0 Flash to add one extra leg to each image in Step 2. We then manually filter these images to retain one high-quality image of animals with clearly either 3 or 5 legs. The final set consists of 91 different animals: 23 three-legged birds and 68 five-legged mammals.

We create three resolution variants for each animal image with size $C \in \{384, 768, 1152\}$ px. Specifically, for an original generated image with dimensions $W \times H$, we resize both dimensions by the *scaling factor* $\frac{C}{\max(W,H)}$ to preserve the original aspect ratio. This procedure generates 91 animals \times 3 resolutions = 273 images in total.

3.2. Tasks 2-5: Counting visual elements in modified familiar patterns: [®] logos, [₱] flags, ² chess pieces, and [®] game boards

Our primary hypothesis is that VLMs contain a strong bias between a brand's logo and its famous visual representations (e.g., an Adidas logo must have 3 stripes; Fig. 2). Expanding beyond animal legs, we test this hypothesis across four domains where humans (and potentially VLMs) have wellestablished visual expectations: logos of famous brands (m), national flag (*), chess pieces (a), and game boards (m). For each domain, we create CF images by making systematic, minimal modifications to familiar visual elements, using the same methodology as Task 1 (m, *) or Python scripts (a, m) with all images rendered at three resolutions (384, 768, and 1152 pixels).

Images For **logos** ^(P) (Appendix D), we modify graphical features (points, prongs, circles, stripes, curves) of three car brands and two shoe brands using ^(D) ○4-mini and ^(D) GPT-4○, placing them in realistic contexts (vehicles and athletic footwear) for a total of 207 images. For **flags** ^(P)

(Appendix E), we systematically add or remove one element (stars or stripes) from 20 flags, creating 120 flag images. For **chess pieces** ² (Appendix F), we test pattern recognition by removing or replacing exactly one piece in western chess and xiangqi starting positions, generating 144 chessboard images. For **game boards** I (Appendix G), we add or remove exactly one row and column across four game types (chess, xiangqi, Sudoku, Go), producing 84 images of the board of these games.

3.3. Task 6: Visual testing on original and modified optical illusion patterns ↔

Recent VLMs show improved performance on optical illusion (Zhang et al., 2023; Guan et al., 2024a) tasks, with IllusionVQA (Shahgir et al., 2024). However, these VLMs might have merely memorized the common 44 optical illusions rather than truly perceiving visual information. To investigate this hypothesis, we create two scenarios: (1) original optical illusions (e.g., the Ebbinghaus illusion where two identical central circles appear different sizes due to surrounding context circles) and (2) modified versions with similar visual setups but reversed effects (e.g., where one central circle is actually larger than the other; Fig. 5). When tested on these modified illusions, VLMs often incorrectly claim the circles are equal (i.e. the answer true for the original illusion but false for the modified version), suggesting a strong bias toward memorized knowledge.

Images We use six classical $\stackrel{\text{Ge}}{\rightarrow}$ optical illusions (Makowski et al., 2021): Müller-Lyer (Müller-Lyer, 1889; Howe & Purves, 2005), Zöllner (Zöllner, 1862; Wallace, 1975), Ebbinghaus (Titchener, 1905; Aglioti et al., 1995), Vertical-Horizontal (Fick, 1851; Hamburger & Hansen, 2010), Pogendorff (Poggendorff, 1863; Green & Hoyle, 1963), and Ponzo (Ponzo, 1910; Yildiz et al., 2022). For five of these illusions, we generate 24 images per type (12 original and 12 modified versions with varying illusion strength). The Vertical-Horizontal illusion, which uses a fixed T-shape that cannot vary in strength, we create only 12 images (6 original and 6 modified). Each image is rendered at three different resolutions: 384, 768, and 1152 pixels. This approach yielded $(24 \times 5 + 12) \times 3 = 396$ images in total. More details in Appendix H.

3.4. Counting the circles or lines in an anomaly cell among a patterned grid **∷**

VLMs can infer the patterns from nearby visual elements to answer visual questions (Huang et al., 2024). We test how much the overall pattern in an image biases its answer to a question about an anomaly region that does *not* obey the pattern.

Images We generate two types of grids-dice and tally

grids—with dimensions $G \times G$, where G ranges from 6 to 12. Cells in dice grids contain circles (Fig. 1#, Fig. F17ab), while tally grids use tally marks (Fig. F17c-d). All grids follow a symmetric visual pattern where the number of shapes in each cell increases from 1 at the edges towards the center, which contains |(G+1)/2| shapes, based on the cell's distance from the nearest edge. For each generated grid image, we introduce an anomaly by modifying **only a** single, strategically chosen cell (avoiding edges & corners). The modifications depend on the grid type: in tally grids, we either add one extra tally mark or remove one existing tally mark from this single anomaly cell; in **dice grids**, we either remove a single circle or replace one circle with a different geometric shape (e.g., triangle, square, star) within this single anomaly cell. Each resulting image is rendered at three resolutions: 384, 768, and 1152 pixels. To create diverse scenarios for single-cell anomaly placement, for each of the 7 grid dimensions, we define two distinct settings by choosing a different single cell to be anomaly, resulting in 14 unique base scenarios (7 dimensions \times 2 choices of a single anomaly cell location).

This systematic generation then yields a total of 2 grid types (dice, tally) \times 2 modification types per grid type \times 14 unique base scenarios for single-cell anomaly placement \times 3 resolutions = 168 distinct images, each featuring one anomaly cell per grid. More details in Appendix I.

4. Results

4.1. Sanity check: VLMs *do* recognize familiar visual subjects on original, unmodified images

Here, we first verify that the subjects in our VLMBias are, in fact, known to VLMs in their original form. If VLMs fail to recognize the concepts in these unaltered images, there is no basis to attribute the their failures on modified images to their bias.

Experiments We evaluate five VLMs (Gemini-2.5 Pro, [Sonnet-3.7, GPT-4.1, o3, and Image: Optimized on the setting of a set of 66 unmodified images spanning our 6 out of 7 VLMBias tasks ([™] animals, [™] logos, [™] flags, ² chess pieces, [™] board game grids, **III** patterned grids). We exclude **III** from the sanity check since the patterns are created from scratch and do not exist on the Internet. For five counting tasks (from ***** to), we ask two questions (identification and counting; Fig. 2a) per image for a total of 132 questions. Since the optical illusion is not a counting task (Fig. $1 \Leftrightarrow$), we replace the counting question with a question asking for the known questions and answers associated with each illusion. That is, we provide an image per illusion type to VLMs (e.g., Fig. 5) and ask VLMs to identify: (1) the name of the illusion; and (2) the question & correct answers associated

Vision Language Models are Biased



Figure 3: VLMs fail to detect subtle changes in counterfactuals (CF) and default to *biased* answers.

with this famous illusion (see the sanity-check prompts in Appendix N).

Results All five VLMs achieve 100% accuracy on all the questions. That is, for counting tasks, VLMs correctly recognize the subjects and the expected counts (e.g., a puma has four legs and the Adidas logo has three stripes; Fig. 3a&d). For all 6 illusion types, VLMs are able to identify the name (e.g., "Ebbinghaus illusion" in Fig. 5), the associated question ("Are the two inner circles equal in size?") and its correct answer ("Yes"). The results here set the ground for the claims in subsequent sections that VLMs' low accuracy on counterfactual images (17.05% accuracy; see Tab. 2) stem from their prior knowledge about the subjects.

4.2. VLMs fail to recognize an extra leg is added to common birds and mammals

Experiments We use the same experiment setup as in Sec. 4.1 but test VLMs on CF images. Specifically, we evaluate five VLMs on the \Im animal images where an extra leg is added to (a) a bird (three legs instead of two) and a mammal (five legs instead of four). We ask each VLM with default settings to count legs (Q1 and Q2; Fig. 2b).

Results On average, VLMs perform poorly (2.12% accuracy) at counting legs of 3-legged and 5-legged counterfactual animals (Tab. 2⁺). Furthermore, 94.14% of the wrong answers match the original, well-known leg counts (Fig. 4 and Fig. 1a), demonstrating that VLMs rely mostly on memorized prior knowledge to answer rather than inspecting the legs in the image (see Fig. 3c).

VLMs are slightly worse at counting the legs of birds than counting the legs of mammals (1.01% vs. 2.50%; Tab. 3^{*}).



Figure 4: On the **counterfactual** images in VLMBias, five VLMs mostly output answers that match the biased choices that we *predefine* for each question, 75.70% of the time, on average.

This biased behavior is the most severe on the leftmost 6 tasks where there are existing prior knowledge on the Internet. Patterned grid \mathbf{iii} is the only task where the visual pattern is created from scratch in this work. Yet, VLMs still are biased 43.45% of the time.

Bird legs (Fig. 1a) are typically thinner, which may make it harder to detect than mammals' legs (Fig. 3b). On birds, except for \bigcirc GPT-4.1, all VLMs score 0% accuracy (Tab. 3).

4.3. VLMs consistently fail to detect subtle changes in familiar subjects [™] **■** [™]

We test whether the biased behavior of VLMs when counting animal legs (Sec. 4.2) also exists in four other domains of man-made subjects: 🐵 logos of famous brands, 🏲 naTable 2: All VLMs achieve 100% on identification and counting tasks with unmodified images, showing that they fully recognize the original version but fail on the counting questions on the modified images (i.e., counterfactuals) in VLMBias. The mean accuracy of five state-of-the-art VLMs on our seven tasks is 17.05%. If 04-mini achieves the highest accuracy (20.25%) which however is still low. VLMs with "thinking" capabilities (If 04-mini, I 03) only slightly outperform non-thinking models (* Gemini-2.5 Pro, Sonnet-3.7, GPT-4.1).

Model	Accur	acy in o	counting	question	s (<mark>Q1</mark> &	2Q2) on	counterfa	ctual images	Unmodified
	a.🎀	b.🝽	c. 🏴	d. 🙎	e. 🖿	f. 🖬	g. III	Task mean	Task mean
♦ Gemini-2.5 Pro	0.00	1.96	10.42	26.74	2.38	49.81	20.83	16.02	100.00
🔯 Sonnet-3.7	0.00	2.72	13.75	9.03	1.79	54.29	34.52	16.59	100.00
GPT−4.1	9.52	9.07	2.50	8.68	0.00	48.61	18.75	13.88	100.00
03	0.92	7.60	5.00	42.71	2.38	50.38	20.54	18.50	100.00
o4−mini	0.18	9.31	14.58	44.10	4.76	51.26	17.56	20.25	100.00
Mean	2.12	6.13	9.25	26.25	2.26	50.87	22.44	17.05	100.00

tional flags, 2 chess pieces, and 2 game boards.

Experiments We replicate the experiments in Sec. 4.2 on CF versions of [™], [™], ², and [™]. For each domain, we create CF images by making systematic modifications: (1) adding or removing a single well-known element (e.g., a stripe in the Adidas logo) in [™] logos; (2) adding or removing a star or a stripe in common [™] national flags; (3) replacing or removing a piece from standard starting chess or xiangqi position; and (4) removing or adding a row or a column in well-known boardgame boards (e.g., Sudoku and Go).

Results VLMs generally demonstrate systematic failures to detect modifications across all four domains, with performance varying depending on the tasks.

- B For logos, accuracy is significantly worse on car logos than on shoe logos (0.44% vs. 17.57%; Tab. 3^(h)). This might be because a logo on a car often appear much smaller than a logo on a shoe photo (Fig. 1^(h) vs. Fig. 2b).
- For flags, VLMs perform better on counting stars (11.79%; Tab. 3⁺) than counting stripes (4.52%; Tab. 3⁺). Counting stripes may be harder because a stripe is often placed right next to other stripes in a flag while stars are spatially separate symbols (Fig. F20b vs. d). More results on flags are in Appendix L.2.
- On counting chess pieces, thinking VLMs (* Gemini-2.5 Pro, and * o4-mini) significantly outperform non-thinking models (>26% vs. <10%; Tab. 32), suggesting that explicit reasoning capabilities help detect anomalies (more results are in Sec. 4.3).
- All VLMs perform extremely poorly (2.26% mean accuracy; Tab. T4) on counting rows and columns of a counterfactual board-game board (Fig. F20c-e). They score 0% accuracy on Sudoku and Go grids (Fig. F13a-b), confirming a fundamental inability to

perform counting (Rahmanzadehgervi et al., 2024), here, in biased counterfactual scenarios (more results in Appendix L.4).

Our findings across four domains collectively show that VLMs rely heavily on memorized knowledge to answer rather than performing detailed visual analysis of the counterfactual image.

4.4. VLMs are biased towards the known illusions and fail to recognize the changes in the counterfactual, modified versions ↔

Experiment We test five VLMs on 6 classic optical illusions, i.e., Müller-Lyer, Zöllner, Ebbinghaus, Vertical-Horizontal, Pogendorff, and Ponzo (Fig. F15). Each illusion is presented in two versions: (a) its original form and (b) a counterfactual, modified version where the groundtruth answer is reversed (Fig. 5). For both versions per illusion, we ask VLMs the same Y/N question (see Appendix H).

Results On average, over original and CF versions, all 5 VLMs perform around the random chance (mean accuracy of 50.87%; Tab. $3 \bowtie$). They tend to provide answers that are *true to the original* versions (i.e., 78.02% mean accuracy) but *false given the counterfactual* versions (23.74% accuracy).

4 out of 5 VLMs perform well on the original versions of the illusions but poorly on the CF versions, exhibiting a strong bias to the well-known illusions. However, Sonnet-3.7 is the opposite—it performs much better on the counterfactual versions than on the original illusions (65.91% vs. 42.68% accuracy; Tab. 3↔). On average, Sonnet-3.7 still performs only slightly above the random chance (54.29% accuracy), revealing a poor low-level vision capability, consistent with recent findings (Rahmanzadehgervi et al., 2024).

Table 3: VLMs perform poorly across	six (out of seven) VLME	Bias tasks, spanning ph	oto-realistic images (🏲 animals and
🝽 logos) and abstract images (🏲 flag,	2 chess pieces, 🛶 optica	al illusions, and 🏬 patter	ned grids).	

		a. 🎀 Animal		b	. 🐨 Log	0		c. 🏴 Flag		d. 🙎 Cł	ess/Xiangq	i Pieces	e. 🚧 (Optical Illusi	ons	f. Hí	Patterned G	rid
Model	Birds	Mammals	Mean	Shoes	Cars	Mean	Stars	Stripes	Mean	Chess	Xiangqi	Mean	Original	Modified	Mean	Remove	Rep/Add	Mean
+	0.00	0.00	0.00	5.80	0.00	1.96	11.54	8.33	10.42	17.36	36.11	26.74	73.16	26.52	49.81	13.10	28.57	20.83
10	0.00	0.00	0.00	8.15	0.00	2.72	20.51	1.19	13.75	7.64	10.42	9.03	42.68	65.91	54.29	35.71	33.33	34.52
99	5.07	11.03	9.52	25.36	1.11	9.07	3.21	1.19	2.50	11.81	5.56	8.68	92.17	5.05	48.61	10.12	27.38	18.75
(0.00	1.23	0.92	21.01	1.11	7.60	5.13	4.76	5.00	56.94	28.47	42.71	91.67	9.09	50.38	14.88	26.19	20.54
8	0.00	0.25	0.18	27.54	0.00	9.31	18.59	7.14	14.58	55.56	32.64	44.10	90.40	12.12	51.26	12.50	22.62	17.56
Mean	1.01	2.50	2.12	17.57	0.44	6.13	11.79	4.52	9.25	29.86	22.64	26.25	78.02	23.74	<mark>50.87</mark>	17.26	27.62	22.44

4.5. VLMs are biased towards the global pattern in a grid **iii**

Experiments We test VLMs on counting the shapes or tally marks inside an anomaly cell where the total number of shapes or marks do not follow the patterns in the surrounding cells (Fig. 1g).

Results Overall, VLMs perform poorly at 22.44% accuracy. 43.45% of all count predictions, both correct and incorrect, match the biased answers (Fig. 411) that correspond to the surrounding cells. In other words, when VLMs make a *wrong* counting predictions, more than half (i.e., 56.02%) of the time, their answers match the **global pattern of most cells in the grid** rather than the target anomaly cell in question (Fig. F17). Our results confirm a striking influence of the background pattern to VLMs' assessment on a small local region. Here, our patterns in the grids are created from scratch and, therefore, do not represent a pattern memorized from the Internet.

4.6. VLMs continue to misidentify the common biased patterns when they do not exist in counterfactual images

Prior sections have shown that VLMs struggle to count the key elements in well-known subjects (e.g., the stripes in a counterfactual, 4-striped Adadias-like logo; Fig. 2b) at 17.05% accuracy (Tab. 2). And \sim 75% of the time, they default to the predefined bias choices. Here, we aim to confirm that VLMs are so biased that they are unable to tell the difference between the original version and the counterfactual by a more direct binary, Yes/No identification question of Q3: "Is this an animal with 4 legs?" when the counterfactual (e.g., a 5-legged puma Fig. 3c) is shown.

Experiments We ask 5 VLMs the Q3 question given our sets of original and CF images. The correct answer is "Yes" for original cases and "No" for all CF cases (Fig. 3c).

Results All VLMs achieve 100% accuracy on the original images, but collapse to a mean of 25.11% on the counterfactual versions (Tab. 4), which is only half of random guessing. That is, VLMs consistently answer "Yes", misidentifying the well-known subject even when the visual evidence con-

Table 4: Accuracy (%) of VLMs on question Q3 (e.g.,, "Is this an animal with 4 legs?") when the image is original (4 legs) or counterfactual (5 legs). VLMs mostly answer "Yes" even on counterfactuals, resulting in accuracy far below the 50% random baseline.

Model	Original	Counterfactual (Δ)						
✦Gemini-2.5 Pro	100.00	20.63 (-79.37)						
🛿 Sonnet-3.7	100.00	23.08 (-76.92)						
₲ GPT-4.1	100.00	26.10 (-73.90)						
Ø 03	100.00	26.15 (-73.85)						
Ø∮04-mini	100.00	29.61 (-70.39)						
Mean	100.00	25.11 (-74.89)						
Answer in curly brackets, e.g., {Yes} or {No}. Original illusion GT: Yes Modified illusion GT: No								
+In-ima	ige text <u>≼</u>	+						
Ebbinghaus illusion	Ebbi	nghaus illusion						
	8							

Figure 5: Original vs. modified versions without (top) and with (bottom) the in-image text ("Ebbinghaus illusion").

tradicts the prompt (Fig. 3c&f). On top of the prior results with Q1 and Q2, the results on Q3 provide extra evidence supporting the hypothesis that VLMs are too biased to recognize that the well-known pattern has changed in counterfactual images.

4.7. Adversarial in-image text showing the name of the common subject further fools VLMs

Prior sections have shown that VLMs perform poorly on the objective task of counting when the background contains **visual** cues strongly correlate with well-known subjects. As VLM outputs may be influenced by adversarial or distracting

text in the image (Goh et al., 2021), here, we test how inimage **textual** cues about the subjects (e.g., "Ebbinghaus illusion") influence VLMs on the same counting questions.

Experiments We insert the subject name (e.g., "Adidas" or "Ebbinghaus illusion"; Fig. 5) into the top of all original and CF images, extending the image vertically but keeping the original content unchanged. We repeat previous experiments asking VLMs the two counting questions (Q1 & Q2).

Results All VLMs perform worse when an in-image text is added (-4.49; Tab. 5). Interestingly, the decrease is more pronounced for thinking models (Tab. 5), such as $4 - \min$ (-6.56), $3 \circ 3$ (-6.41), than for non-thinking ones such as Sonnet-3.7 (-2.81) and GPT-4.1 (-2.67). This result is consistent with recent findings that thinking models tend to hallucinate more (OpenAI, 2025a; Zhang et al., 2024a), here more biased toward the text in the image despite contradictory visuals.

Table 5: Adding adversarial, in-image textual cues that state the subject name (e.g., "Adidas") cause VLMs to decrease their accuracy (-4.49) on counterfactual images (b). In contrast, instructing VLMs to rely exclusively on the image details to answer questions (Debiased) or to double-check its answers (Double-Check) only slightly improves accuracy, by +1.87 and +2.70, respectively (c).

Model	a. Baseline	b. Adversarial	c. Helpful textual prompt		
		w/ In-image text	w/ Debiased Prompt	w/ Double-Check	
♦ Gemini-2.5 Pro	16.02	12.04 (-3.98)	19.72 (+3.70)	20.22 (+4.20)	
🛯 Sonnet-3.7	16.59	13.78 (-2.81)	19.29 (+2.70)	20.86 (+4.27)	
₿ GPT-4.1	13.88	11.21 (-2.67)	14.38 (+0.50)	16.00 (+2.12)	
🙆 o 3	18.50	12.09 (-6.41)	18.94 (+0.44)	21.02 (+2.52)	
Ø∱o4-mini	20.25	13.69 (-6.56)	22.25 (+2.00)	20.61 (+0.36)	
Mean	17.05	12.56 (-4.49)	18.92 (+1.87)	19.75 (+2.70)	

4.8. Helpful prompts do not ameliorate the bias issues in VLM

Previous results show that VLMs rely heavily on prior knowledge to answer objective counting questions. Here, we test how incorporating *helpful* instructions in the prompts may help VLMs become less biased.

Experiments We apply two prompting strategies across all VLMBias tasks:

(1) Debiased Prompt: We prepend the original question (Q1 and Q2) with "*Do not assume from prior knowledge and answer only based on what is visible in the image.*" to encourage models to rely exclusively on image contents.

(2) Double-Check: After VLMs answer the original question, we add a follow-up prompt of "*Please double-check your answer and give your final answer in curly brackets, following the format above.*"

These prompts are designed to encourage VLMs to examine the image more carefully. All experiments use the same images and default model settings as in the baseline setup.

Results Both helpful prompting strategies improve VLM accuracy but only slightly over the baseline, +1.87 for Debiased and +2.70 for Double-Check (Tab. 5c). That is, explicitly instructing models to rely on image contents or verify their answer helps to some extent but does not address the core issue of bias.

5. Discussion and Conclusion

Our study shows that SOTA VLMs fail consistently in counting visual elements (e.g., stripes in a logo) when they are strongly biased towards the subject (e.g., an Adidas logo has three stripes). Thinking models (0, 04-mini, 0, 03) perform only slightly better than nonthinking ones (0 Sonnet-3.7, \clubsuit Gemini-2.5 Pro, 0 GPT-4.1)— longer thoughts do not address the bias issue. Similarly, instructing VLMs to double check its answers or rely exclusively on the image contents only modestly increase the accuracy. Our work documents important visual biases of VLMs in an objective counting task rather than the common social biases (Sheng et al., 2019a) often documented in the literature.

For 6 out of 7 tasks, we use an automated pipeline consisting of scripts, LLMs and text-to-image generators that generate counterfactual images. In our pipeline, humans do not manually edit the original images to create counterfactuals but only review the generated images. We release the automation scripts and evaluation code.

It might be interesting to compare whether the original and counterfactual images map to similar visual representations after the vision encoders in VLMs. That is, are the VLM failures in this paper the result of visual encoders unable to capture the fine-grained modifications in the couterfactual images? Alternatively, vision encoders may be able to observe the visual changes but the LLMs in late-fusion architectures are too biased to output accurate answers.

Limitations Our work has two limitations. First, VLMs with image generation capabilities (e.g., and) are still in early developmental stages and exhibit *their own biases*, making it non-trivial to control generated images as expected. This limitation prevents us from exploring some other interesting domains. Second, we are unable to test VLMs that have the capability to use tools, which may substantially help on VLMBias questions. That is, it might be interesting to test how the biased visual cues in VLMBias images suppress tool-use VLMs to use their tools (e.g., zooming functions) to answer questions. Our preliminary results with the chat interface of o3 reveal that o3 often does not even use its visual thinking capability (OpenAI, 2025b) to examine the counterfactual images but instead directly attempt to answer questions.

Acknowledgement

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT)(RS-2025-00573160), and Innovative Human Resource Development for Local Intellectualization program through the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT)(IITP-2025-RS-2020-II201489).

We also thank Khang Gia Le (Independent Researcher) for feedback and discussions of the earlier results. AV was supported by Hyundai Motor Chung Mong-Koo Global Scholarship. AN was supported by the NSF Grant No. 1850117 & 2145767, and donations from NaphCare Foundation & Adobe Research.

References

- Abid, A., Farooqi, M., and Zou, J. Persistent anti-muslim bias in large language models. In Fourcade, M., Kuipers, B., Lazar, S., and Mulligan, D. K. (eds.), AIES '21: AAAI/ACM Conference on AI, Ethics, and Society, Virtual Event, USA, May 19-21, 2021, pp. 298–306. ACM, 2021. doi: 10.1145/3461702.3462624. URL https: //doi.org/10.1145/3461702.3462624.
- Aglioti, S., DeSouza, J. F., and Goodale, M. A. Size-contrast illusions deceive the eye but not the hand. *Current biology*, 5(6):679–685, 1995.
- Anonymous. B-score: Detecting biases in large language models using response history. In Fortysecond International Conference on Machine Learning, 2025. URL https://openreview.net/forum? id=kl7SbPfBsB.
- Anthropic. Claude 3.7 Sonnet and Claude Code, 2025. URL https://www.anthropic.com/ news/claude-3-7-sonnet. https://www. anthropic.com/news/claude-3-7-sonnet.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-VL: A versatile visionlanguage model for understanding, localization, text reading, and beyond. arXiv preprint arXiv:2308.12966, 2023.
- Bang, Y., Chen, D., Lee, N., and Fung, P. Measuring political bias in large language models: What is said and how it is said. In Ku, L., Martins, A., and Srikumar, V. (eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pp. 11142–11159. Association for Computational Linguistics, 2024. doi: 10.18653/V1/ 2024.ACL-LONG.600. URL https://doi.org/ 10.18653/v1/2024.acl-long.600.
- Bitton-Guetta, N., Bitton, Y., Hessel, J., Schmidt, L., Elovici, Y., Stanovsky, G., and Schwartz, R. Breaking common sense: Whoops! a vision-and-language benchmark of synthetic and compositional images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2616–2627, 2023.
- Echterhoff, J. M., Liu, Y., Alessa, A., McAuley, J., and He, Z. Cognitive bias in decision-making with LLMs. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 12640–12653, Miami, Florida, USA, November 2024. Association for Computational Linguistics. URL https://aclanthology.org/ 2024.findings-emnlp.739.

- Fick, A. De errone quodam optic asymmetria bulbi effecto. *Marburg: Koch*, 1851.
- Fraser, K. and Kiritchenko, S. Examining gender and racial bias in large vision-language models using a novel dataset of parallel images. In Graham, Y. and Purver, M. (eds.), Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 690–713, St. Julian's, Malta, March 2024. Association for Computational Linguistics. URL https://aclanthology. org/2024.eacl-long.41/.
- Gallegos, I. O., Rossi, R. A., Barrow, J., Tanjim, M. M., Kim, S., Dernoncourt, F., Yu, T., Zhang, R., and Ahmed, N. K. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3):1097–1179, 2024.
- Goh, G., Cammarata, N., Voss, C., Carter, S., Petrov, M., Schubert, L., Radford, A., and Olah, C. Multimodal neurons in artificial neural networks. *Distill*, 6(3):e30, 2021.
- Google. Google Gemini 2.5 Pro, 2025. URL https://deepmind.google/technologies/ gemini/pro/. https://deepmind.google/ technologies/gemini/pro/.
- Green, R. and Hoyle, E. The poggendorff illusion as a constancy phenomenon. *Nature*, 200(4906):611–612, 1963.
- Guan, T., Liu, F., Wu, X., Xian, R., Li, Z., Liu, X., Wang, X., Chen, L., Huang, F., Yacoob, Y., et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 14375– 14385, 2024a.
- Guan, T., Liu, F., Wu, X., Xian, R., Li, Z., Liu, X., Wang, X., Chen, L., Huang, F., Yacoob, Y., et al. HallusionBench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *CPVR*, 2024b.
- Hall, S. M., Gonçalves Abrantes, F., Zhu, H., Sodunke, G., Shtedritski, A., and Kirk, H. R. Visogender: A dataset for benchmarking gender bias in image-text pronoun resolution. Advances in Neural Information Processing Systems, 36:63687–63723, 2023.
- Hamburger, K. and Hansen, T. Analysis of individual variations in the classical horizontal-vertical illusion. *Attention*, *Perception*, & *Psychophysics*, 72(4):1045–1052, 2010.

- Hirota, Y., Nakashima, Y., and Garcia, N. Gender and racial bias in visual question answering datasets. In *Proceedings* of the 2022 ACM Conference on Fairness, Accountability, and Transparency, pp. 1280–1292, 2022.
- Howard, P., Madasu, A., Le, T., Moreno, G. L., Bhiwandiwalla, A., and Lal, V. Socialcounterfactuals: Probing and mitigating intersectional social biases in vision-language models with counterfactual examples. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11975–11985, 2024.
- Howe, C. Q. and Purves, D. The müller-lyer illusion explained by the statistics of image–source relationships. *Proceedings of the National Academy of Sciences*, 102 (4):1234–1239, 2005.
- Hu, T., Kyrychenko, Y., Rathje, S., Collier, N., van der Linden, S., and Roozenbeek, J. Generative language models exhibit social identity biases. *Nat. Comput. Sci.*, 5(1):65–75, 2025. doi: 10.1038/ S43588-024-00741-1. URL https://doi.org/10. 1038/s43588-024-00741-1.
- Huang, W., Liu, H., Guo, M., and Gong, N. Visual hallucinations of multi-modal large language models. In Ku, L., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pp. 9614–9631. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024. FINDINGS-ACL.573. URL https://doi.org/10.18653/v1/2024.findings-acl.573.
- Janghorbani, S. and De Melo, G. Multi-modal bias: Introducing a framework for stereotypical bias assessment beyond gender and race in vision-language models. In Vlachos, A. and Augenstein, I. (eds.), Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pp. 1725–1735, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.eacl-main. 126. URL https://aclanthology.org/2023. eacl-main.126/.
- Koevering, K. V. and Kleinberg, J. M. How random is random? evaluating the randomness and humaness of llms' coin flips. *CoRR*, abs/2406.00092, 2024. doi: 10.48550/ARXIV.2406.00092. URL https://doi. org/10.48550/arXiv.2406.00092.
- Koo, R., Lee, M., Raheja, V., Park, J. I., Kim, Z. M., and Kang, D. Benchmarking cognitive biases in large language models as evaluators. In Ku, L., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and*

virtual meeting, August 11-16, 2024, pp. 517–545. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-ACL.29. URL https://doi.org/10.18653/v1/2024.findings-acl.29.

- Kumar, D., Jain, U., Agarwal, S., and Harshangi, P. Investigating implicit bias in large language models: A large-scale study of over 50 llms, 2024. URL https://arxiv.org/abs/2410.12864.
- Lee, K.-i., Kim, M., Yoon, S., Kim, M., Lee, D., Koh, H., and Jung, K. VLind-Bench: Measuring language priors in large vision-language models. In *NAACL Findings*, 2025.
- Lee, N., Bang, Y., Lovenia, H., Cahyawijaya, S., Dai, W., and Fung, P. Survey of social bias in vision-language models. *arXiv preprint arXiv:2309.14381*, 2023.
- Li, C., Chen, M., Wang, J., Sitaram, S., and Xie, X. Culturellm: Incorporating cultural differences into large language models. *Advances in Neural Information Processing Systems*, 37:84799–84838, 2024.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. Advances in neural information processing systems, 36:34892–34916, 2023.
- Liu, J., Fu, Y., Xie, R., Xie, R., Sun, X., Lian, F., Kang, Z., and Li, X. Phd: A chatgpt-prompted visual hallucination evaluation dataset. *arXiv preprint arXiv:2403.11116*, 2024.
- Makowski, D., Lau, Z. J., Pham, T., Boyce, W. P., and Chen, S. A. A parametric framework to generate visual illusions using python. *Perception*, 50(11):950–965, 2021. doi: 10.1177/03010066211057347. URL https://doi. org/10.1177/03010066211057347. PMID: 34841973.
- Müller-Lyer, F. C. Optische Urteilstäuschungen. Archiv für Anatomie und Physiologie, Physiologische Abteilung, 2: 263–270, 1889. Original description of the Müller-Lyer illusion.
- Naous, T., Ryan, M. J., Ritter, A., and Xu, W. Having beer after prayer? measuring cultural bias in large language models. In Ku, L., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 16366–16393. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.862. URL https://doi.org/10.18653/v1/2024.acl-long.862.
- OpenAI. Introducing GPT-4.1 in the API, 2025a. URL https://openai.com/index/gpt-4-1/. https://openai.com/index/gpt-4-1/.

OpenAI. Introducing OpenAI o3 and o4-mini, 2025b. URL https://openai.com/ index/introducing-o3-and-o4-mini/. https://openai.com/index/ introducing-o3-and-o4-mini/.

- OpenAI. Openai o3 and o4-mini system card. Technical report, OpenAI, April 2025a. URL https://cdn.openai.com/pdf/ 2221c875-02dc-4789-800b-e7758f3722c1/ o3-and-o4-mini-system-card.pdf. Comprehensive system card detailing the capabilities, safety, and evaluation results for OpenAI o3 and o4-mini models.
- OpenAI. Thinking with images, 2025b. URL https://openai.com/index/ thinking-with-images/. Accessed: 2025-05-28.
- Parrish, A., Chen, A., Nangia, N., Padmakumar, V., Phang, J., Thompson, J., Htut, P. M., and Bowman, S. R. BBQ: A hand-built bias benchmark for question answering. In Muresan, S., Nakov, P., and Villavicencio, A. (eds.), *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 2086–2105. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022. FINDINGS-ACL.165. URL https://doi.org/10.18653/v1/2022.findings-acl.165.
- Poggendorff, J. C. Biographisch-literarisches handwörterbuch zur geschichte der exakten wissenschaften von jc poggendorff, i-ii. Leipzig: Johann Ambrosius Barth. Ponatis:(1965). Amsterdam: BM Israël NV, 1863.
- Ponzo, M. Intorno ad alcune illusioni nel campo delle sensazioni tattili, sull'illusione di Aristotele e fenomeni analoghi. Wilhelm Engelmann, 1910.
- Potter, Y., Lai, S., Kim, J., Evans, J., and Song, D. Hidden persuaders: Llms' political leaning and their influence on voters. In Al-Onaizan, Y., Bansal, M., and Chen, Y. (eds.), Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024, pp. 4244–4275. Association for Computational Linguistics, 2024. URL https://aclanthology.org/ 2024.emnlp-main.244.
- Rahmanzadehgervi, P., Bolton, L., Taesiri, M. R., and Nguyen, A. T. Vision language models are blind. In Cho, M., Laptev, I., Tran, D., Yao, A., and Zha, H. (eds.), Computer Vision - ACCV 2024 - 17th Asian Conference on Computer Vision, Hanoi, Vietnam, December 8-12, 2024, Proceedings, Part V, volume 15476 of Lecture Notes in Computer Science, pp. 293–309. Springer, 2024.

doi: 10.1007/978-981-96-0917-8_17. URL https://doi.org/10.1007/978-981-96-0917-8_17.

- Raj, C., Mukherjee, A., Caliskan, A., Anastasopoulos, A., and Zhu, Z. Biasdora: Exploring hidden biased associations in vision-language models. In Al-Onaizan, Y., Bansal, M., and Chen, Y. (eds.), *Findings of the Association for Computational Linguistics: EMNLP* 2024, Miami, Florida, USA, November 12-16, 2024, pp. 10439–10455. Association for Computational Linguistics, 2024. URL https://aclanthology.org/2024. findings-emnlp.611.
- Ruggeri, G. and Nozza, D. A multi-dimensional study on bias in vision-language models. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 6445–6455, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023. findings-acl.403. URL https://aclanthology. org/2023.findings-acl.403/.
- Sathe, A., Jain, P., and Sitaram, S. A unified framework and dataset for assessing societal bias in vision-language models. In Al-Onaizan, Y., Bansal, M., and Chen, Y. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pp. 1208–1249. Association for Computational Linguistics, 2024. URL https://aclanthology. org/2024.findings-emnlp.66.
- Shahgir, H. S., Sayeed, K. S., Bhattacharjee, A., Ahmad, W. U., Dong, Y., and Shahriyar, R. Illusionvqa: A challenging optical illusion dataset for vision language models. arXiv preprint arXiv:2403.15952, 2024.
- Sharma, P., Shaham, T. R., Baradad, M., Fu, S., Rodriguez-Munoz, A., Duggal, S., Isola, P., and Torralba, A. A vision check-up for language models. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14410–14419, 2024.
- Sheng, E., Chang, K., Natarajan, P., and Peng, N. The woman worked as a babysitter: On biases in language generation. In Inui, K., Jiang, J., Ng, V., and Wan, X. (eds.), Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019, pp. 3405–3410. Association for Computational Linguistics, 2019a. doi: 10.18653/V1/D19-1339. URL https://doi.org/10.18653/v1/D19-1339.
- Sheng, E., Chang, K.-W., Natarajan, P., and Peng, N. The woman worked as a babysitter: On biases in language generation. In Inui, K., Jiang, J., Ng, V., and Wan, X.

(eds.), Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 3407–3412, Hong Kong, China, November 2019b. Association for Computational Linguistics. doi: 10.18653/v1/D19-1339. URL https://aclanthology.org/D19-1339/.

- Shin, J., Song, H., Lee, H., Jeong, S., and Park, J. Ask llms directly, "what shapes your bias?": Measuring social bias in large language models. In Ku, L., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pp. 16122–16143. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-ACL. 954. URL https://doi.org/10.18653/v1/ 2024.findings-acl.954.
- Taesiri, M. R., Nguyen, G., Habchi, S., Bezemer, C., and Nguyen, A. Imagenet-hard: The hardest images remaining from a study of the power of zoom and spatial biases in image classification. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023. URL http://papers. nips.cc/paper_files/paper/2023/hash/ 706390d6f9208b03bc54f97ac3cfe99e-Abstract and_Benchmarks.html.
- Team, C. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- Titchener, E. B. *Experimental psychology: A manual of laboratory practice*, volume 2. Macmillan Company, 1905.
- Tong, S., Liu, Z., Zhai, Y., Ma, Y., LeCun, Y., and Xie, S. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *CVPR*, 2024.
- Wallace, G. The effect of contrast on the zöllner illusion. *Vision Research*, 15(8-9):963–966, 1975.
- Wang, W., Jiao, W., Huang, J., Dai, R., Huang, J., Tu, Z., and Lyu, M. R. Not all countries celebrate thanksgiving: On the cultural dominance in large language models. In Ku, L., Martins, A., and Srikumar, V. (eds.), *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pp. 6349–6384. Association for Computational Linguistics, 2024. doi: 10. 18653/V1/2024.ACL-LONG.345. URL https://doi. org/10.18653/v1/2024.acl-long.345.

- Wu, X., Guan, T., Li, D., Huang, S., Liu, X., Wang, X., Xian, R., Shrivastava, A., Huang, F., Boyd-Graber, J. L., Zhou, T., and Manocha, D. Autohallusion: Automatic generation of hallucination benchmarks for vision-language models. In Al-Onaizan, Y., Bansal, M., and Chen, Y. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pp. 8395–8419. Association for Computational Linguistics, 2024. URL https://aclanthology.org/2024.findings-emnlp.493.
- Xiao, Y., Liu, A., Cheng, Q., Yin, Z., Liang, S., Li, J., Shao, J., Liu, X., and Tao, D. Genderbias-vl: Benchmarking gender bias in vision language models via counterfactual probing. *CoRR*, abs/2407.00600, 2024. doi: 10.48550/ ARXIV.2407.00600. URL https://doi.org/10. 48550/arXiv.2407.00600.
- Yildiz, G. Y., Sperandio, I., Kettle, C., and Chouinard, P. A. A review on various explanations of ponzo-like illusions. *Psychonomic Bulletin & Review*, pp. 1–28, 2022.
- Zhang, Y., Pan, J., Zhou, Y., Pan, R., and Chai, J. Grounding visual illusions in language: Do vision-language models perceive illusions like humans? In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 5718–5728. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023. EMNLP-MAIN.348. URL https://doi.org/10.18653/v1/2023.emnlp-main.348.
- Zhang, Y., Li, Y., Wang, Y., Wang, X., Wang, Y., and Wang, X. How language model hallucinations can snowball. In Proceedings of the 41st International Conference on Machine Learning, volume 235 of PMLR, pp. 59670–59684, 2024a. URL https://proceedings.mlr.press/v235/zhang24ay.html. Shows how step-by-step reasoning can propagate and amplify hallucinations in large language models.
- Zhang, Y., Schwarzschild, A., Carlini, N., Kolter, Z., and Ippolito, D. Forcing diffuse distributions out of language models. *CoRR*, abs/2404.10859, 2024b. doi: 10.48550/ ARXIV.2404.10859. URL https://doi.org/10. 48550/arXiv.2404.10859.
- Zhao, J., Fang, M., Pan, S., Yin, W., and Pechenizkiy, M. GPTBIAS: A comprehensive framework for evaluating bias in large language models. *CoRR*, abs/2312.06315, 2023. doi: 10.48550/ARXIV.2312.06315. URL https://doi.org/10.48550/arXiv.2312.06315.
- Zhou, K., Lai, E., Yeong, W. B. A., Mouratidis, K., and Jiang, J. ROME: Evaluating pre-trained vision-language

models on reasoning beyond visual common sense. In *Findings of the Association for Computational Linguistics: EMNLP*, 2023. URL https://openreview. net/forum?id=N6sXsHuWDE.

Zöllner, F. Ueber eine neue art anorthoskopischer zerrbilder. *Annalen der Physik*, 193(11):477–484, 1862. doi: https://doi.org/10.1002/andp.18621931108. URL https://onlinelibrary.wiley.com/doi/ abs/10.1002/andp.18621931108.

Appendix for: Vision Language Models are Biased: Counting legs of an animal is surprisingly hard

Table of Contents

• Illustrative questions
Models and access details
Task 1: Counting legs with added limb
Task 2: Counting elements in modified brand logos
Task 3: Counting stripes/stars in modified national flags
• Task 4: Counting chess pieces on modified starting position
Task 5: Counting rows and columns of board game
• Task 6: Visual testing with both original and modified optical illusion
• Task 7: Counting circles or lines in an anomaly cell within a patterned grid
• Qualitative results on 🕆 animals
• Qualitative results on F flags
• More findings
Prompts used for image generation and image editing
Questions for sanity check

A. Illustrative questions

Торіс	Subtopic	Q1	Q2	Q3
Animal		How many legs does this animal have? Answer with a number in curly brackets, e.g., {9}.	Count the legs of this ani- mal. Answer with a num- ber in curly brackets, e.g., {9}.	Is this an animal with 4 legs? Answer in curly brackets, e.g., {Yes} or {No}.
Logo	Adidas	How many visible stripes are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}.	Count the visible stripes in the logo on the left shoe. Answer with a number in curly brackets, e.g., {9}.	Are the logos on these shoes Adidas logos? An- swer in curly brackets, e.g., {Yes} or {No}.
	Nike	How many visible white stylized curves are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}.	Count the visible white stylized curves in the logo on the left shoe. Answer with a number in curly brackets, e.g., {9}.	Are the logos on these shoes Nike logos? An- swer in curly brackets, e.g., {Yes} or {No}.
	Mercedes	How many points are there on the star in the logo of this car? Answer with a number in curly brackets, e.g., {9}.	Count the points on the star in the logo of this car. Answer with a number in curly brackets, e.g., {9}.	Is the logo on this car Mercedes-Benz logo? An- swer in curly brackets, e.g., {Yes} or {No}.
	Audi	How many overlapping circles are there in the logo of this car? Answer with a number in curly brackets, e.g., {9}.	Count the overlapping circles in the logo of this car. Answer with a number in curly brackets, e.g., {9}.	Is the logo on this car Audi logo? Answer in curly brackets, e.g., {Yes} or {No}.
	Maserati	How many prongs are there in the logo of this car? Answer with a num- ber in curly brackets, e.g., {9}.	Count the prongs in the logo of this car. Answer with a number in curly brackets, e.g., {9}.	Is the logo on this car Maserati logo? Answer in curly brackets, e.g., {Yes} or {No}.
Flag	Stars	How many stars are there on this flag? Answer with a number in curly brackets, e.g., {9}.	Count the stars on this flag. Answer with a number in curly brackets, e.g., {9}.	Is this the flag of [coun- try]? Answer in curly brackets, e.g., {Yes} or {No}.
	Stripes	How many stripes are there on this flag? Answer with a number in curly brackets, e.g., {9}.	Count the stripes on this flag. Answer with a number in curly brackets, e.g., {9}.	Is this the flag of [coun- try]? Answer in curly brackets, e.g., {Yes} or {No}.

Table T1: Some examples of questions on 🏲 animal, 🐨 brand logos, and 🏲 flags

Торіс	Subtopic	Q1	Q2	Q3
Chess Pieces	Chess	How many chess pieces are there on this board? Answer with a number in curly brackets, e.g., {9}.	Count the chess pieces on this board. Answer with a number in curly brackets, e.g., {9}.	Is this the chess starting position? Answer in curly brackets, e.g., {Yes} or {No}.
	Xiangqi	How many xiangqi pieces are there on this board? Answer with a number in curly brackets, e.g., {9}.	Count the xiangqi pieces on this board. Answer with a number in curly brackets, e.g., {9}.	Is this the Xiangqi start- ing position? Answer in curly brackets, e.g., {Yes} or {No}.
Board Game	Chess	How many rows are there on this board? Answer with a number in curly brackets, e.g., {9}.	Count the rows on this board. Answer with a number in curly brackets, e.g., {9}.	Is this a 8x8 Chessboard? Answer in curly brackets, e.g., {Yes} or {No}.
	Xiangqi	How many horizontal lines are there on this board? Answer with a number in curly brackets, e.g., {9}.	Count the horizontal lines on this board. Answer with a number in curly brackets, e.g., {9}.	Is this a 10x9 Xiangqi board? Answer in curly brackets, e.g., {Yes} or {No}.
	Go	How many horizontal lines are there on this board? Answer with a number in curly brackets, e.g., {9}.	Count the horizontal lines on this board. Answer with a number in curly brackets, e.g., {9}.	Is this a 19x19 Go board? Answer in curly brackets, e.g., {Yes} or {No}.
	Sudoku	How many rows are there on this puzzle? Answer with a number in curly brackets, e.g., {9}.	Count the rows on this puz- zle. Answer with a num- ber in curly brackets, e.g., {9}.	Is this a 9x9 Sudoku puz- zle? Answer in curly brackets, e.g., {Yes} or {No}.
Patterned Grid	Dice	How many circles are there in cell C5? Answer with a number in curly brackets, e.g., {9}.	Count the circles in cell C5. Answer with a number in curly brackets, e.g., {9}.	Does cell C5 contain 4 circles? Answer in curly brackets, e.g., {Yes} or {No}.
	Tally	How many lines are there in cell C5? Answer with a number in curly brackets, e.g., {9}.	Count the lines in cell C5. Answer with a number in curly brackets, e.g., {9}.	Does cell C5 contain 3 lines? Answer in curly brackets, e.g., {Yes} or {No}.

Table T2: Some examples of questions on 2 chesse pieces, 🕮 game boards and 🗰 patterned grid.

Торіс	Subtopic	Q1	Q2	Q3
Optical Illusion	Ebbinghaus	Are the two inner circles equal in size? Answer in curly brackets, e.g., {Yes} or {No}.	Do the two inner circles have the same size? An- swer in curly brackets, e.g., {Yes} or {No}.	Is this an example of the Ebbinghaus illusion? Answer in curly brackets, e.g., {Yes} or {No}.
	Mullerlyer	Are the two horizontal lines equal in length? An- swer in curly brackets, e.g., {Yes} or {No}.	Do the two horizontal lines have the same length? Answer in curly brackets, e.g., {Yes} or {No}.	Is this an example of the Müller-Lyer illusion? An- swer with Yes/No. An- swer in curly brackets, e.g., {Yes} or {No}.
	Poggendorff	Are the two diagonal line segments aligned? An- swer in curly brackets, e.g., {Yes} or {No}.	Do the two diagonal lines form a straight line? An- swer in curly brackets, e.g., {Yes} or {No}.	Is this an example of the Poggendorff illusion? An- swer in curly brackets, e.g., {Yes} or {No}.
	Ponzo	Are the two horizontal lines equal in length? An- swer in curly brackets, e.g., {Yes} or {No}.	Do the two horizontal lines have the same length? Answer in curly brackets, e.g., {Yes} or {No}.	Is this an example of the Ponzo illusion? Answer in curly brackets, e.g., {Yes} or {No}.
	VerticalHorizontal	Are the horizontal and ver- tical lines equal in length? Answer in curly brackets, e.g., {Yes} or {No}.	Do the horizontal and ver- tical lines have the same length? Answer in curly brackets, e.g., {Yes} or {No}.	Is this an example of the Vertical–Horizontal il- lusion? Answer in curly brackets, e.g., {Yes} or {No}.
	Zollner	Are the two horizontal lines parallel? Answer in curly brackets, e.g., {Yes} or {No}.	Do the two horizontal lines run parallel? An- swer in curly brackets, e.g., {Yes} or {No}.	Is this an example of the Zöllner illusion? Answer in curly brackets, e.g., {Yes} or {No}.

Table T3: Some examples of questions on $\stackrel{\smile}{\rightarrow}$ optical illusions.

B. Models and access details

We evaluate five state-of-the-art VLMs using their official APIs with default settings. These include three thinking models (e.g., Gemini-2.5 Pro, 03, 04-mini) and two non-thinking models (e.g., Sonnet-3.7, GPT-4.1).

B.1. Gemini-2.5 Pro

We access Gemini-2.5 Pro (gemini-2.5-pro-preview-05-06) via aistudio.google.com and use all *default* settings with temperature=1.0.

B.2. Sonnet-3.7

We access the Anthropic API via console.anthropic.com to use Sonnet-3.7 (claude-3.7-sonnet) and *default* settings with temperature=1.0.

B.3. GPT-4.1

We access the API for GPT-4.1 (gpt-4.1) via platform.openai.com and use all *default* settings with temperature=1.0.

B.4. o3

We access the OpenAI API for o3 via platform.openai.com and use *default* settings with temperature=1.0.

B.5. 04-mini

We access the OpenAI API for o4-mini (o4-mini) via platform.openai.com with *default* settings including:

- temperature: 1.0
- reasoning_effort: medium (default thinking mode setting)

C. Task 1: Counting legs with added limb >>

C.1. Task design



Figure F1: Data generation pipeline for Task 1: Counting legs with added limb.

Pretrained on the Internet data, VLMs must have colossal prior knowledge of the count of animal legs from both textual and image data. Following this hypothesis, we generate images of usual animals with *one additional leg* (e.g., 3-legged birds or 5-legged dogs) and ask VLMs to count legs to evaluate if these models are biased toward their prior knowledge.

- Animal types: We modify the legs of 2 types of animals: birds and mammals.
- Modification types: Each animal is modified to have 1 additional leg.
- Target animals: We select 91 well-known animals, consisting of 23 two-legged birds and 68 four-legged mammals.
- **Image resolutions**: We generate each animal image and rescale them at **3** different pixel sizes {384, 768, 1152}px using the scaling factor in Sec. **3**.1 to test resolution sensitivity

This approach generates a total of 91 animals \times 1 modification type \times 3 resolutions = 273 total images.

C.2. Implementation and image generation

Implementation details Our image generation pipeline follows this sequence:

- 1. Use 1. Use -mini to collect a list of well-known animals with clearly visible legs
- 2. Generate full-body and side-view images of these animals using +/Gemini-2.0 Flash
- 3. For each animal image, use +/Gemini-2.0 Flash to add one extra leg to the animal. Each animal image is edited over 4 independent trials.
- 4. Manually inspect and filter out unsatisfactory images
- 5. Render each approved image at three different resolutions

Quality control We manually inspect the images to ensure that each modified animal image has exactly one additional leg. For cases that fail (e.g., more than one added leg), we remove them from our dataset.

Prompt We use the following prompts to test the VLMs:

- Q1: How many legs does this animal have? Answer with a number in curly brackets, e.g., {9}.
- Q2: Count the legs of this animal. Answer with a number in curly brackets, e.g., {9}.
- Q3: Is this an animal with [NumModifiedLegs] legs? Answer in curly brackets, e.g., {Yes} or {No}.

Ground truth calculation The ground truth answers are as follow:

- Birds leg counting (Q1&Q2):
 - Correct answer: 3 (one additional leg)
 - Expected bias: 2
- Mammals leg counting (Q1&Q2):
 - Correct answer: 5 (one additional leg)
 - Expected bias: 4
- Animal leg identification question (Q3):
 - Correct answer: "No" (always, since each animal has one additional leg)
 - Expected bias: "Yes"

C.3. Qualitative results



Figure F2: VLMs are often biased toward the original number of legs \Im animals have, and they tend to answer based on prior knowledge rather than by analyzing the image.

C.4. List of animals

Mammals: Four-legged animals

horse, zebra, donkey, mule, cow, buffalo, yak, water buffalo, deer, elk, moose, reindeer, caribou, gazelle, giraffe, camel, dromedary camel, bactrian camel, llama, alpaca, goat, ibex, mountain goat, pronghorn, bighorn sheep, wild boar, pig, warthog, coyote, lynx, bobcat, leopard, tiger, lion, jaguar, puma, ocelot, caracal, hyena, rabbit, impala, springbok, kudu, eland, wildebeest, okapi, hippopotamus, african elephant, asian elephant, indian rhinoceros, gnu, maned wolf, arctic fox, red fox, fennec fox, red wolf, domestic dog, domestic cat, african wilddog, dingo, jackal, gray wolf, hare, cheetah, antelope, bison, sheep, serval

Birds: Two-legged animals

ostrich, emu, rhea, cassowary, heron, stork, crane, egret, ibis, spoonbill, turkey, chicken, rooster, duck, swan, peacock, sandpiper, avocet, stilt, plover, lapwing, oystercatcher, secretary bird



D. Task 2: Counting elements in modified brand logos 🕫

Figure F3: Data generation pipeline of shoe logos for Task 2: Counting elements in modified brand logos



Figure F4: Data generation pipeline of car logos for Task 2: Counting elements in modified brand logos

D.1. Task design

Our initial evaluation show that some VLMs, such as $3^{\circ} \circ 4$ -mini, can accurately count the four stripes on modified Adidas logo on white background. As such, to increase the task difficulty, we hypothesize that VLMs strongly associate 3° logos with the background they typically appear on. Subsequently, we examine if the visual cues from the background would be strong enough to suppress counting the elements in the logos. Our task is designed as follow:

- Brand types: We use 2 different brand types: cars and shoes
- Target brands: We select 5 well-known brands with quantifiable graphical elements:
 - Car brands: Mercedes-Benz, Maserati, and Audi (3 brands)

Shoe logos

(b): How many visible stripes are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}. (b): How many visible white stylized curves are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}.

($\widehat{\mathbf{w}}$ (d): How many visible black stylized curves are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., $\{9\}$.

	(a) Adidas		(b) Nike	(0	c) Adidas		(d) Nike	(e) Adidas
					3		e e		8	
+	3	×	1	×	3	×	2	 Image: A second s	3	×
*	3	×	1	×	3	×	1	×	3	×
\$	3	X	1	×	3	×	1	×	3	×
9	3	×	1	×	3	×	1	×	4	 Image: A second s
6	3	×	1	×	3	×	1	×	4	\checkmark
Bias	3	×	1	×	3	×	1	×	3	×
GT	4	√	2	 Image: A second s	4	√	2	√	4	√
	4	Gemini-2	.5 Pr	o 📲 Sonne	et-3.7	₲ GPT-4	4.1 🧕	o3 🚳 🗲	04-mini	-

Figure F5: VLMs are often biased and rely on prior knowledge when answering questions about [®] shoe logos, even with simple ones like the Nike Swoosh. Please zoom in to see the logo clearly.

- Shoe brands: Adidas and Nike (2 brands)
- **Background variations**: Each brand logo has specific background settings:
 - Car logo background: Car logos always appear on cars. For each logo, we collect 5 car body types × 3 colors (white, grey, black)
 - Shoe logo background: Shoe logos are often seen on the footwear of athletes. For each logo, we collect a list of 4 relevant sports (tennis, running, basketball, soccer) × 3 colors (black, red, white)
- **Image resolutions**: We generate each image and rescale them at **3** different pixel sizes {384, 768, 1152}px using the scaling factor in Sec. 3.1 to test resolution sensitivity

This systematic approach generates a total of [3 car brands \times (5 \times 3) \times 3 resolutions] + [2 shoe brands \times (4 \times 3) \times 3 resolutions] = 135 + 72 = 207 total images.

D.2. Implementation and prompts

Implementation details We employ the following process to generate logo modification images:

- 1. Use Use !-. Use !-. Use !-. Use !-. We then select the most relevant suggestions for our benchmark.
- 2. Generate modified logo versions using GPT-40.
- 3. Create background images:
 - Background images for car logos:

- Use @fo4-mini to suggest popular colors and body types of each car logo.
- For each logo, generate and select relevant images of cars from the logo brand with the determined body types and colors.
- Manually place modified logos in typical car logo positions.
- Background images for shoe logos:

 - For each logo, generate and select relevant images of athletes wearing shoes with the modified logo for each determined color and sport.
- 4. Render each image at three different resolutions.

Quality control To ensure high-quality images, we manually review to make sure that: (1) each generated logo has the correct number of modified elements; (2) each product is clearly visible and oriented correctly; and (3) the position of the logos on the products are natural-looking.

Prompts We use the following prompts

- 1. Counting questions (Q1 & Q2):
 - Q1 (Adidas): *How many visible stripes are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}.*
 - Q1 (Nike): How many visible [CurveColor] stylized curves are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}
 - Q1 (Audi): How many overlapping circles are there in the logo of this car? Answer with a number in curly brackets, e.g., {9}.
 - Q1 (Mercedes): How many points are there on the star in the logo of this car? Answer with a number in curly brackets, e.g., {9}.
 - Q1 (Maserati): How many prongs are there in the logo of this car? Answer with a number in curly brackets, e.g., {9}
 - **Q2** (Adidas): Count the visible stripes in the logo of the left shoe. Answer with a number in curly brackets, e.g., *{9}.*
 - Q2 (Nike): Count the visible [CurveColor] stylized curves in the logo of the left shoe. Answer with a number in curly brackets, e.g., {9}
 - Q2 (Audi): Count the overlapping circles in the logo of this car. Answer with a number in curly brackets, e.g., {9}.
 - Q2 (Mercedes): Count the points on the star in the logo of this car. Answer with a number in curly brackets, e.g., {9}.
 - Q2 (Maserati): Count the prongs in the logo of this car. Answer with a number in curly brackets, e.g., {9}
- 2. Y/N identification questions (Q3):
 - Q3 (Adidas): Are the logos on these shoes Adidas logos? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q3 (Nike): Are the logos on these shoes Nike logos? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q3 (Audi): Is the logo on this car Audi logo? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q3 (Mercedes): Is the logo on this car Mercedes-Benz logo? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q3 (Maserati): Is the logo on this car Maserati logo? Answer in curly brackets, e.g., {Yes} or {No}.

Ground truth calculation The ground truth answers are as follow:

- Adidas stripes counting (Q1&Q2):
 - Correct answer: 4
 - Expected bias: 3
- Nike stylized curves counting (Q1&Q2):

- Correct answer: 2
- Expected bias: 1

• Audi overlapping circles counting (Q1&Q2):

- Correct answer: 5
- Expected bias: 4
- Mercedes-Benz points on the star counting (Q1&Q2):
 - Correct answer: 4
 - Expected bias: 3
- Maserati prongs counting (Q1&Q2):
 - Correct answer: 5
 - Expected bias: 3
- Logo identification question (Q3):
 - Correct answer: "No" (all logos are modified)
 - Expected bias: "Yes"

4

4

4

4

5

Х

Х

Х

X

3

3

3

3

4

✦ Gemini-2.5 Pro

х

х

X

X

1

\$

6

Bias

GT

D.3. Qualitative results

Car logos (a), (d): How many overlapping circles are there in the logo of this car? Answer with a number in curly brackets, e.g., {9}. (b), (c): How many points are there on the star in the logo of this car? Answer with a number in curly brackets, e.g., {9}. 🐨 c: How many prongs are there in the logo of this car? Answer with a number in curly brackets, e.g., {9}. (b) Mercedes (e) Mercedes (a) Audi (c) Maserati (d) Audi 3 3 4 4 X 3 柴 4 Х 3 Х 3 Х 4 X 3

3

3

3

3

5

Sonnet-3.7

4

4

4

4

5

X

Х

X

X

o3

3

3

3

3

4

х

Х

Х

Х

Х

X

X

1

₲ GPT-4.1



Shoe logos

(a), (c), (e): How many visible stripes are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}. (b): How many visible white stylized curves are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., {9}.

m (d): How many visible black stylized curves are there in the logo of the left shoe? Answer with a number in curly brackets, e.g., $\{9\}$.



Figure F7: VLMs are often biased and rely on prior knowledge when answering questions about [®] shoe logos, even with simple ones like the Nike Swoosh. Please zoom in to see the logo clearly.



Figure F8: Data generation pipeline for Task 3: Counting stripes/stars in modified national flags.

E. Task 3: Counting stripes/stars in modified national flags 🏴

E.1. Task design

Flags of countries contain easily recognizable patterns. To evaluate if existing VLMs overly rely on their knowledge of these **F** flags to count a certain element, we design the task as follow:

- Flag types: We modify 2 commonly used elements across different flags: stars and stripes
- Modification types: Each flag has 2 types of modifications:
 - Add: We add an additional element (star or stripe) to a chosen flag
 - Remove: We remove one element (star or stripe) from a chosen flag
- **Target flags**: We select 20 well-known country flags with either 3+ stars or 5+ stripes (a total of 13 star-typed flags and 7 stripe-typed flags) to ensure the modified flags retain recognizable traits to test visual bias.
- Image resolutions: We generate each flag and rescale them at 3 different pixel sizes {384, 768, 1152}px using the scaling factor in Sec. 3.1 to test resolution sensitivity

This systematic approach generates a total of 20 target flags \times 2 modification types \times 3 resolutions = 120 total images.

E.2. Implementation and image generation

Implementation details We modify the SVG code of a chosen flag to create new variants following this sequence:

- 1. Identify 20 well-known country flags (13 with 3+ stars, 7 with 5+ stripes) based on the suggestions from store -mini.
- 2. Retrieve original SVG code from WikiCommons for each flag.
- 3. Use 5-04-mini to modify each SVG to create two variants:
 - An "Add" variant with one additional element.

- A "Remove" variant with one fewer element.
- 4. Render each modified flag at three different resolutions.

Quality control We employ the following steps to ensure high-quality and consistent images:

- Manual inspection: We manually review each generated sample to verify modification quality and visual consistency
- Filtering: We remove unsatisfactory samples from the benchmark and rerun the pipeline on these cases to obtain new samples.
- **Fallback**: For rare cases (3 in total) that consistently fail automated generation, we manual modify the flags to ensure they strictly follow the modification rules.

Prompts We use the following prompts:

- 1. Counting questions (Q1 & Q2):
 - Q1 (Star-typed flags): How many stars are there on this flag? Answer with a number in curly brackets, e.g., {9}.
 - Q1 (Stripe-typed flags): *How many stripes are there on this flag? Answer with a number in curly brackets, e.g., {9}.*
 - Q2 (Star-typed flags): Count the stars on this flag. Answer with a number in curly brackets, e.g., {9}.
 - Q2 (Stripe-typed flags): Count the stripes on this flag. Answer with a number in curly brackets, e.g., {9}.
- 2. Y/N identification questions (Q3):
 - Is this the flag of [CountryName]? Answer in curly brackets, e.g., {Yes} or {No}.

Ground truth calculation We calculate the ground truth as follow:

- Direct counting questions (Q1 & Q2):
 - Correct answer: The actual count of the elements (stars or stripes) on the flag after modification
 - * For *Remove modifications*: Standard element count minus 1
 - * For Add modifications: Standard element count plus 1
 - Expected bias: The standard element count
- Flag verification question (Q3):
 - Correct answer: "No" (since the flag's element has been modified)
 - Expected bias: "Yes"

National Flag



Figure F9: VLMs are biased when counting the stars and stripes on **F** national flags.



F. Task 4: Counting chess pieces on modified starting position a

Figure F10: Data generation pipeline for Task 4: Counting chess pieces on modified starting position

F.1. Task design

To evaluate if VLMs rely on expected structure or attend to actual pieces, we test their ability to count pieces on boards with subtle modifications. We design our task with careful control of visual parameters to ensure systematic evaluation:

- Board types: We use 2 different game boards: { chess (Western chess), xiangqi (Chinese chess) }
- Modification types: Each board has 2 types of modifications:
 - Remove: We remove exactly one piece from the standard starting position.
 - Replace: We replace exactly one piece with a different piece of the same color.
- **Target squares**: We select 12 unique occupied squares per board type, maintaining the same target squares across the Remove and Replace modifications to ensure controlled comparison.
- Image resolutions: We generate each board at 3 different pixel sizes {384, 768, 1152}px to test resolution sensitivity.

This systematic approach generates a total of 2 board types \times 2 modification types \times 12 target squares \times 3 resolutions = 144 total images.

F.2. Implementation and prompts

Implementation details Our implementation utilizes specialized libraries for each board type. For chess, we leverage the Python chess library to manipulate board states and chess.svg for rendering. For xiangqi (Chinese chess), we created a custom implementation using svgwrite for rendering.

The algorithm for both board types follows the same sequence:

- 1. Create a standard board with all 32 pieces in their starting positions
- 2. Randomly select 12 target squares from the occupied squares
- 3. For each target square, create (a) a Remove variant and (b) a Replace variant
- 4. Render each modified board at three different resolutions

The xiangqi implementation required special handling for:

• The traditional 9×10 board layout with the central river and two palaces

- · Chinese character rendering for pieces, which requires detecting appropriate CJK fonts
- Different piece distribution (Chariots, Knights, Elephants, Advisors, General, Cannons, and Soldiers)

Quality control To ensure consistent image quality across all variants, we implemente several technical measures:

- SVG to PNG conversion: We used direct SVG rendering with adjustable scaling factors based on target resolution
- Quality scaling: We applied a quality multiplier $(5.0 \times \text{ base resolution factor})$ to ensure clear piece visibility

Prompts We use different prompts for each modification type to test VLMs' visual attention:

1. Remove modifications:

- Q1: How many [chess/xiangqi] pieces are there on this board? Answer with a number in curly brackets, e.g., {9}.
- Q2: Count the [chess/xiangqi] pieces on this board. Answer with a number in curly brackets, e.g., {9}.

2. Replace modifications:

- **Q1:** *How many* [Added Piece Type] pieces are there on this board? Answer with a number in curly brackets, e.g., {9}.
- Q2: Count the [Added Piece Type] pieces on this board? Answer with a number in curly brackets, e.g., {9}.

3. Both modification types:

• Q3: Is this the [chess/xiangqi] starting position? Answer in curly brackets, e.g., {Yes} or {No}.

For Replace modifications, [Added Piece Type] refers to the specific piece type that is added to the board through replacement, chosen from:

- For chess: Pawn, Knight, Bishop, Rook, Queen, or King
- For xiangqi: Soldier, Horse, Elephant, Chariot, Cannon, Advisor, or General

For Replace modifications, we ask about the added piece type rather than total count because this more effectively tests whether VLMs rely on prior knowledge of standard piece distributions or actually inspect the board carefully.

Ground truth calculation We calculate the ground truth answers for each prompt type:

- Total piece count (Remove modifications only):
 - Correct answer: 31 (one fewer than the standard 32 pieces)
 - Expected bias: 32 (the standard piece count)
- Added piece type count (Replace modifications only):
 - Correct answer: The standard count for that piece type plus one
 - For example, if a Knight is replaced with a Bishop in chess, the Bishop count would be 3 (standard 2 + 1 added)
 - Expected bias: The standard count for that piece type (e.g., 2 for Bishops in chess)
 - This tests if VLMs rely on their knowledge of standard piece counts or actually inspect the board
- Starting position question (Both modification types):
 - Correct answer: Always "No" (since the board has been modified)
 - Expected bias: "Yes" (since the board closely resembles the starting position)

F.3. Qualitative results



Figure F11: VLMs are biased when counting the pieces on 2 chess and xiangqi.



G. Task 5: Counting rows and columns of board game

Figure F12: Data generation pipeline for Task 5: Counting rows and columns of board game

G.1. Task design

To evaluate VLMs' over-reliance on visual bias versus actual counting, we adapted the row and column counting task from BlindTest (Rahmanzadehgervi et al., 2024) where Claude-3.5-Sonnet achieved 74.26% accuracy. Instead of simple grids, we leverage modified versions of well-known game boards to test whether VLMs rely on prior knowledge or perform actual visual counting. We design our task with careful control of visual parameters to ensure systematic evaluation:

- **Board types**: We use 4 different grid-based game boards: {*Chess* (8×8), *Xiangqi* (Chinese chess, 10×9), *Sudoku* (9×9), *Go* (19×19)}
- Modification types: Each board has up to 4 types of modifications:
 - Remove row: We remove exactly one row from the grid.
 - *Remove column*: We remove exactly one column from the grid.
 - Add row: We add exactly one row to the grid.
 - Add column: We add exactly one column to the grid.
- **Board-specific variations**: For Chess, Xiangqi, and Sudoku boards, all four modifications (remove/add row, remove/add column) are visually distinct, with additional positional variations (first/last), resulting in 8 variants per board. Go boards have uniform grid structure, so we produce only 4 variations.
- Image resolutions: We generate each board at 3 different pixel sizes {384, 768, 1152}px to test resolution sensitivity.

This systematic approach generates a total of (8 variants \times 3 board types (Xiangqi/Chess/Sudoku) + 4 Go variants) \times 3 resolutions = 84 total images.

G.2. Implementation and prompts

Implementation details Our implementation utilizes specialized drawing libraries for each board type. For Chess, we use standard 8×8 chessboard grid generation with alternating square colors. For Xiangqi, we implement the traditional 10×9

board layout with river gap and palace diagonal lines. For Sudoku, we create 9×9 grids with bold 3×3 block boundaries and sample numbers. For Go, we generate uniform line grids with traditional star points.

The algorithm for all board types follows the same sequence:

- 1. Create a standard board with correct dimensions and visual elements
- 2. Apply systematic modifications (add/remove rows/columns at specific positions)
- 3. Maintain visual consistency of special elements
- 4. Render each modified board at three different resolutions

The board-specific implementations required special handling for:

- · Chess: Alternating light/dark square pattern preservation across dimension changes
- Xiangqi: River gap positioning and palace diagonal lines adjustment for row modifications
- Sudoku: Bold 3×3 block boundary lines based on original 9×9 grid structure
- Go: Uniform line spacing and star point positioning for various board sizes

Quality control To ensure consistent image quality across all variants, we implemented several technical measures:

- SVG to PNG conversion: We used direct SVG rendering with adjustable scaling factors based on target resolution
- Quality scaling: We applied a quality multiplier (5.0× base resolution factor) to ensure clear structural visibility
- Font and layout fidelity: Automatic detection and usage of appropriate fonts, particularly critical for Xiangqi (Chinese characters) and Sudoku (numbers)

Table T4: All VLMs' performance is extremely low (2.26%) across \blacksquare game boards, confirming that current VLMs are largely unable to perform even simple counting operations in structured visual settings

Model	Chess	Go	Sudoku	Xiangqi	Mean
✦Gemini-2.5 Pro	2.08	0.00	0.00	6.25	2.38
Sonnet-3.7	0.00	0.00	0.00	6.25	1.79
₲ GPT-4.1	0.00	0.00	0.00	0.00	0.00
◎ o3	0.00	0.00	0.00	8.33	2.38
	16.67	0.00	0.00	0.00	4.76
Mean	3.75	0.00	0.00	4.17	2.26

Prompts We use different prompts for different question types to test VLMs' visual counting versus prior knowledge:

1. Counting questions (Q1 & Q2):

- Q1 (Chess): *How many* [rows/columns] are there on this board? Answer with a number in curly brackets, e.g., {9}.
- Q1 (Xiangqi, Go): *How many [horizontal/vertical] are there on this board? Answer with a number in curly brackets, e.g., {9}.*
- Q1 (Sudoku): *How many* [rows/columns] are there on this puzzle? Answer with a number in curly brackets, e.g., {9}.

- Q2 (Chess): Count the [rows/columns] on this board. Answer with a number in curly brackets, e.g., {9}.
- Q2 (Xiangqi, Go): Count the [horizontal/vertical] lines on this board. Answer with a number in curly brackets, e.g., {9}.
- Q2 (Sudoku): Count the [rows/columns] on this puzzle. Answer with a number in curly brackets, e.g., {9}.

2. Y/N identification questions (Q3):

- Q3 (Chess): Is this a 8×8 Chessboard? Answer in curly brackets, e.g., {Yes} or {No}.
- Q3 (Xiangqi): Is this a 10×9 Xiangqi board? Answer in curly brackets, e.g., {Yes} or {No}.
- Q3 (Sudoku): Is this a 9×9 Sudoku puzzle? Answer in curly brackets, e.g., {Yes} or {No}.
- Q3 (Go): Is this a 19×19 Go board? Answer in curly brackets, e.g., {Yes} or {No}.

Ground truth calculation We calculate the ground truth answers for each prompt type:

- Row/Column count (Q1 & Q2):
 - Correct answer: The actual number of rows/columns after modification. For example, if one row is removed from a 9×9 Sudoku, the row count is 8.
 - **Expected bias**: The standard count for that board type (e.g., 8 for Chess rows, 10 for Xiangqi horizontal lines, 9 for Sudoku rows, 19 for Go horizontal lines)
- Standard layout question (Q3):
 - Correct answer: Always "No" (since all boards have been modified from standard dimensions)
 - Expected bias: "Yes" (since the boards closely resemble their standard counterparts)

Game Boards

(a): How many columns are there on this puzzle? Answer with a number in curly brackets, e.g., {9}.
(b), (c): How many horizontal lines are there on this board? Answer with a number in curly brackets, e.g., {9}.
(d): How many rows are there on this board? Answer with a number in curly brackets, e.g., {9}.



Figure F13: VLMs are biased when counting the rows and columns on \mathbb{H} game boards.



H. Task 6: Visual testing with both original and modified optical illusion

Figure F14: Data generation pipeline for Task 6: Visual testing with both original and modified optical illusion

H.1. Task design

Recent VLMs show improved performance on optical illusion tasks, with o4-mini achieving 71.49% accuracy on IllusionVQA. However, these VLMs might have merely memorized the common optical illusions rather than truly perceiving visual information. To investigate this hypothesis, we test their ability to correctly identify illusion effects on both original and strategically modified versions. We design our task with careful control of visual parameters to ensure systematic evaluation:

- **Illusion types**: We use **6** different classical optical illusions: {*Ebbinghaus*, *Müller-Lyer*, *Ponzo*, *Vertical-Horizontal*, *Zöllner*, *Poggendorff* }
- Condition types: Each illusion has 2 conditions:
 - *Original*: Standard illusion where the visual effect should occur (e.g., two identical circles appearing different sizes).
 - *Modified*: Reversed version where the actual measurements contradict the typical illusion effect (e.g., circles that are genuinely different sizes).
- Parameter variations: We generate multiple combinations of illusion parameters:
 - Most illusions: 12 original + 12 modified versions with varying illusion strength and difference
 - Vertical-Horizontal: 6 original + 6 modified versions (fixed T-shape structure)
- Image resolutions: We generate each illusion at 3 different pixel sizes {384, 768, 1152}px to test resolution sensitivity.

This systematic approach generates a total of (12 original + 12 modified) \times 5 illusion types + (6 original + 6 modified) \times 1 Vertical-Horizontal illusion) \times 3 resolutions = 396 total images.

H.2. Implementation and prompts

Implementation details Our implementation adapts code from Pyllusion (https://github.com/ RealityBending/Pyllusion) to generate consistent, parametrically controlled optical illusions. We systematically vary two key parameters: *illusion strength* (which controls the intensity of contextual elements that create the illusion effect, representing how strongly the surrounding context biases perceptual experience) and *difference* (which controls the objective, actual difference between target elements being compared, where 0 means identical elements and non-zero values create genuine physical differences).

The algorithm for all illusion types follows the same sequence:

1. Define parameter ranges for each illusion type (strength values, difference values).

- 2. Generate original versions with standard illusion parameters (diff=0 for equal elements).
- 3. Generate modified versions with reversed parameters (diff $\neq 0$ for unequal elements).
- 4. Render each illusion variant at three different resolutions.

The illusion-specific implementations required special parameter handling for:

- Ebbinghaus: Varying surrounding circle sizes (strength) and central circle differences (difference).
- Müller-Lyer: Different arrowhead angles (strength) and line length differences (difference).
- Ponzo: Perspective line angles (strength) and horizontal bar length differences (difference).
- Vertical-Horizontal: Fixed T-shape with varying line length ratios (difference).
- Zöllner: Background line angles (strength) and main line parallelism differences (difference).
- Poggendorff: Interrupting rectangle positions (strength) and diagonal line alignments (difference).

Quality control To ensure consistent image quality and valid illusion effects across all variants, we implemented several technical measures:

- **Parameter validation**: Ensured all strength and difference values produce visually meaningful illusions, with diff ≠ 0 cases design to be easily recognizable by humans to distinguish actual physical differences from perceptual biases clearly.
- Balanced generation: Equal numbers of diff=0 (original) and diff \neq 0 (modified) cases per illusion type

Prompts We use consistent prompts across illusion types to test VLMs' visual perception versus memorized knowledge:

- 1. Main questions (Q1 & Q2):
 - Q1 (Ebbinghaus): Are the two inner circles equal in size? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q1 (Müller-Lyer, Ponzo): Are the two horizontal lines equal in length? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q1 (Vertical-Horizontal): Are the horizontal and vertical lines equal in length? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q1 (Zöllner): Are the two horizontal lines parallel? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q1 (Poggendorff): Are the two diagonal line segments aligned? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q2 (Ebbinghaus): Do the two inner circles have the same size? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q2 (Müller-Lyer): Do the two horizontal lines have the same length? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q2 (Ponzo): Do the two horizontal lines have the same length? Answer in curly brackets, e.g., {Yes} or {No}.
 - **Q2** (Vertical-Horizontal): Do the horizontal and vertical lines have the same length? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q2 (Zöllner): Do the two horizontal lines run parallel? Answer in curly brackets, e.g., {Yes} or {No}.
 - Q2 (Poggendorff): Do the two diagonal lines form a straight line? Answer in curly brackets, e.g., {Yes} or {No}.
- 2. Y/N identification questions (Q3):
 - Q3: Is this an example of the [Ebbinghaus/Müller-Lyer/Ponzo/Vertical-Horizontal/Zöllner/Poggendorff] illusion? Answer in curly brackets, e.g., {Yes} or {No}.

Ground truth calculation We calculate the ground truth answers based on the actual measurements in each image:

• Counting questions (Q1 & Q2):

- Correct answer:
 - * Original illusions (diff=0): Elements are actually equal, so the correct answer is "Yes"
 - * Modified illusions (diff $\neq 0$): Elements are actually different, so the correct answer is "No"
- Expected bias:
 - * **Original illusions**: VLMs might incorrectly say "No" expecting the illusion effect to make equal elements appear different
 - * **Modified illusions**: VLMs might incorrectly say "Yes" expecting the illusion to make genuinely different elements appear equal
- Y/N identification questions (Q3):
 - Correct answer:
 - * Original illusions: "Yes" (standard examples of the specified illusion type).
 - * Modified illusions: "No" (modified versions that contradict typical illusion effects).
 - Expected bias:
 - * Original illusions: VLMs likely correctly identify as "Yes" since they match memorized illusion patterns
 - * **Modified illusions**: VLMs may incorrectly say "Yes" if they rely on visual similarity rather than recognizing the effect contradiction

H.3. Qualitative results

Abstract images: Optical Illusions												
	(a) Or Müller	iginal r-Lyer	(b) Mo Mülle	o dified r-Lyer	lified (c) Original Lyer Zöllner		(d) Modified Zöllner		(e) Original Ebbinghaus		(f) Modified Ebbinghaus	
	\succ	\prec	\succ	\prec	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	<i>~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~</i>	~~~~	~~~				
	\leftarrow	\rightarrow	\leftarrow	\rightarrow								•
+	Yes	√	Yes	×	Yes	√	Yes	×	Yes	√	Yes	×
*	Yes	1	Yes	×	Yes		Yes	×	No	×	No	
\$	Yes	√	Yes	×	Yes	√	Yes	×	Yes	 Image: A second s	Yes	×
9	Yes		Yes	×	Yes		Yes	×	Yes		Yes	×
\$ <mark>;</mark>	Yes	1	Yes	×	Yes	 Image: A second s	Yes	×	No	×	Yes	×
Bias	No	×	Yes	×	No	×	Yes	×	No	×	Yes	×
GT	Yes	1	No	1	Yes	 Image: A second s	No	 Image: A second s	Yes	\checkmark	No	√
 ◆ Gemini-2.5 Pro Sonnet-3.7 General GPT-4.1 General of OPT-4.1 Genera												

Figure F15: VLMs show systematic biases, often relying on prior knowledge about \ominus optical illusions rather than directly interpreting the image.



I. Task 7: Counting circles or lines in an anomaly cell within a patterned grid **#**

Figure F16: Data generation pipeline for Task 7: Counting circles or lines in an anomaly cell within a patterned grid

I.1. Task design

VLMs can infer patterns from nearby visual elements to answer visual questions (Huang et al., 2024). To evaluate whether VLMs rely on pattern recognition over actual visual counting, we create square grids with systematic numerical patterns (represented visually by dice faces or tally marks) where exactly one cell violates the expected pattern. We hypothesize that VLMs will prioritize the inferred pattern over the actual visual information and report the expected pattern-completing value instead of the true count. We design our task with careful control of visual parameters to ensure systematic evaluation:

- **Grid types**: We use 2 different visual representation types: {*dice* (circular dots in dice-face patterns), *tally* (traditional tally mark lines)}.
- Modification types per grid type: For each grid type, we apply 2 distinct types of cell-level modifications:
 - *Dice grids*: Remove (one dot is removed from a cell) and Replace (one dot is replaced with a different shape, like a square or star, within a cell).
 - Tally grids: Remove (one tally line is removed from a cell) and Add (one extra tally line is added to a cell).
- Grid Dimensions: We generate grids of 7 different dimensions, ranging from 6×6 to 12×12 cells.
- Unique scenarios for anomaly placement (single anomaly per grid image): To create 14 distinct base settings for placing anomalies, where each final grid image will feature only a single modified cell. We proceed as follows: for each of the 7 grid dimensions, we define two separate base settings. Each of these two settings for a given grid dimension involves selecting a *different*, unique cell location to be the *sole* anomaly cell for images generated under that specific setting. These potential anomaly cell locations are carefully chosen to avoid edges and corners. This gives us (7 grid dimensions \times 2 distinct choices of a single anomaly cell location per dimension) = 14 distinct base settings. For each of these 14 base settings (defined by a grid dimension and the location of its single anomaly cell), we then apply all combinations of grid types and their respective modifications to generate the final images, each still containing only that one pre-determined anomaly.
- **Image resolutions**: Each generated grid image is rendered at **3** different pixel sizes {384, 768, 1152}px to assess sensitivity to image resolution.

This systematic generation process yields a total of 2 (grid types) \times 2 (modification types) \times 14 (unique scenarios) \times 3 (resolutions) = 168 distinct images.

I.2. Implementation and prompts

Implementation details Our implementation generates systematic pattern grids using a distance-from-edge algorithm to create naturally increasing-then-decreasing numerical patterns. For dice grids, we use circular dots arranged in traditional dice-face configurations (1-6 dots per cell). For tally grids, we render authentic tally marks with proper grouping (four vertical lines crossed by a diagonal fifth line).

The algorithm for both grid types follows the same sequence:

- 1. Generate base grid with pattern-consistent cell counts using distance-from-edge calculation
- 2. Organize target positions across 14 groups, with each group containing both dice and tally variants
- 3. For each target cell, create modification variants:
 - Dice: Remove one dot OR replace one dot with alternative shape (triangle, square, star)
 - Tally: Remove one line OR add one extra line
- 4. Render each modified grid at three different resolutions with consistent visual quality

The grid-specific implementations required special handling for:

- Dice pattern consistency: Maintaining standard dice-face arrangements (1-6 dots) while allowing single-dot modifications
- Tally mark authenticity: Proper grouping of marks with diagonal crosses for every fifth line
- Pattern calculation: Distance-from-edge algorithm ensuring natural numerical progression across grid cells
- Cell positioning: Strategic selection of anomaly cells away from edges to preserve pattern context

Quality control To ensure consistent image quality and valid pattern recognition challenges across all variants, we implemented several technical measures:

- SVG to PNG conversion: We used direct SVG rendering with adjustable scaling factors based on target resolution
- Quality scaling: We applied a quality multiplier (5.0× base resolution factor) to ensure clear shape and line visibility

Prompts We use consistent prompts across both grid types to test VLMs' pattern recognition versus actual visual counting:

1. Counting questions (Q1 & Q2):

- Q1 (Dice): How many circles are there in cell [CellID]? Answer with a number in curly brackets, e.g., {9}.
- Q1 (Tally): How many lines are there in cell [CellID]? Answer with a number in curly brackets, e.g., {9}.
- Q2 (Dice): Count the circles in cell [CellID]. Answer with a number in curly brackets, e.g., {9}.
- Q2 (Tally): Count the lines in cell [CellID]. Answer with a number in curly brackets, e.g., {9}.

2. Y/N identification questions (Q3):

- Q3 (Dice): Does cell [CellID] contain [ExpectedCount] circles? Answer in curly brackets, e.g., {Yes} or {No}.
- Q3 (Tally): Does cell [CellID] contain [ExpectedCount] lines? Answer in curly brackets, e.g., {Yes} or {No}.

For all prompts, [CellID] refers to the specific anomaly cell using standard spreadsheet notation (e.g., C3, F7), and [ExpectedCount] represents the pattern-consistent count that would be expected based on surrounding cells.

Ground truth calculation We calculate the ground truth answers based on the actual visual content in each modified cell:

- Direct counting questions (Q1 & Q2):
 - Correct answer: The actual count of visual elements in the target cell after modification
 - * For Remove modifications: Standard pattern count minus 1
 - * For *Add modifications*: Standard pattern count plus 1
 - * For Replace modifications: Standard pattern count minus 1 (since one circle is replaced with a different shape)
 - Expected bias: The pattern-consistent count that VLMs might infer from surrounding cells, ignoring the actual modification
- Pattern-based verification question (Q3):
 - Correct answer: Always "No" (since the target cell has been modified to break the pattern)
 - Expected bias: "Yes" (if VLMs rely on pattern inference rather than direct visual inspection)

I.3. Qualitative results



Figure F17: All VLMs, except 🖾 Sonnet-3.7, fail to correctly identify the abnormal cell (C3) in both the 🗰 patterned grids.

J. Qualitative results on 🐂 animals

How many legs does this animal have? Answer with a number in curly brackets, e.g., {9}.

Baseline: How many legs does this animal have? Answer with a number in curly brackets, e.g., {9}. **w/ Double-Check**: Please double-check your answer and give your final answer in curly brackets, following the format above.

w/ Debiased Prompt: Do not assume from prior knowledge and answer only based on what is visible in the image. Ground truth: $3\checkmark$



← Gemini-2.5 Pro(Baseline): {2} 🗡

```
◆ Gemini-2.5 Pro(w/ Double-Check): {2} ★
◆ Gemini-2.5 Pro(w/ Debiased Prompt): {2} ★
```

Sonnet-3.7 (Baseline): The image shows a rooster (male chicken) standing on a grassy field. Chickens have two legs, which are clearly visible in this image. The rooster is standing upright on its two pale yellow/tan legs with its feet planted on the grass. $\{2\} \times$

Sonnet-3.7 (w/ Double-Check): Looking at the image again, I can clearly see the rooster is standing on two legs. Chickens are bipedal animals, with two legs that they use for walking and standing. In this particular photo, both of the chicken's legs are visible beneath its body. $\{2\} \times$

Sonnet-3.7 (w/ Debiased Prompt): Based solely on what is visible in the image, this chicken has {2} legs. The image clearly shows a rooster standing on a grassy field with its two yellow/tan-colored legs and feet visible beneath its body.

```
    GPT-4.1 (Baseline): {4} ×
    GPT-4.1 (w/ Double-Check): {4} ×
    GPT-4.1 (w/ Debiased Prompt): {4} ×
```

③ o3 (Baseline): {2} ×
③ o3 (w/ Double-Check): {2} ×
③ o3 (w/ Debiased Prompt): {4} ×

Figure F18: VLMs fail 100% of the time, even on simple tasks like counting chicken legs, despite helpful prompts such as debiasing or double-checking.

K. Qualitative results on **F** flags

How many stripes are there in this flag? Answer with a number in curly brackets, e.g., {9}.

Baseline: How many stripes are there in this flag? Answer with a number in curly brackets, e.g., {9}. **w/ Double-Check**: Please double-check your answer and give your final answer in curly brackets, following the format above.

w/ Debiased Prompt: Do not assume from prior knowledge and answer only based on what is visible in the image. Ground truth: $12 \checkmark$



Figure F19: VLMs fail and remain biased toward 13 even when one stripe is removed from the U.S. flag.

L. More findings



Figure F20: VLMs perform poorly at **counting** elements on counterfactual images across ₱, , and a domains, heavily defaulting to the biased answers.

L.1. VLMs struggle to detect logo modifications, often relying on context rather than visual detail 🐵

Experiments We replicate the experiment from Sec. 4.2 on our [®] logo task, evaluating five VLMs on modified shoe and car logo images.

Results VLM performance on car logos $(0.44\%; \text{Tab. 3}^{(1)})$ is significantly worse than on shoe logos $(17.57\%; \text{Tab. 3}^{(2)})$, as the emblem is small relative to the vehicle (see Fig. 1b). In contrast, shoe logos occupy more image area (see Fig. 3e) and involve only a few simple curves or stripes (i.e., one extra curve for Nike, one added stripe for Adidas). These results highlight two key limitations: VLMs fail to attend to small, context-embedded visual changes and instead rely on memorization, without visually verifying the B logo itself (e.g., by zooming in (Taesiri et al., 2023)).

L.2. VLMs fail to count visual elements in modified flags 🏴

Experiments We follow the procedure from Sec. 4.2 on our **F** flag tasks. Five VLMs are prompted to count either the number of stars or the number of stripes in original and modified versions of national flags. Modifications consist of adding or removing a single star or stripe, and each model uses its default settings.

Results VLMs achieve higher mean accuracy on star modifications (11.79%; Tab. $3^{\texttt{H}}$) than on stripe modifications (4.52%; Tab. $3^{\texttt{H}}$). This pattern indicates that models are somewhat more attuned to discrete symbol changes (missing or extra stars; see Fig. F20d) than to subtle structural alterations (extra or missing stripes; see Fig. F20b), yet overall sensitivity to flag modifications is extremely limited (9.25%; Tab. $3^{\texttt{H}}$).

L.3. Thinking models better detect piece changes in modified chess starting positions 2

Experiments We evaluate five VLMs on a ^a chess-piece counting task using standard starting positions for both Western chess and xiangqi. For each board type, we generate images in which exactly one piece is either removed or replaced by

another piece of the same color. All models use their default settings and are prompted to report the total number of pieces or number of a certain piece (e.g., Knights) on the board.

Results VLMs perform significantly better on Western chess (see Fig. 1²) than on xiangqi (see Fig. F20a) in terms of mean accuracy (29.86 % vs. 22.64%; Tab. 3²). Thinking models (Gemini-2.5 Pro, 30, and 4-mini) all exceed 26% accuracy, whereas non-thinking models (GPT-4.1 and Sonnet-3.7) remain below 10% (Tab. 3²). This suggests that on well-structured abstract images, models with explicit reasoning capabilities are better able to detect anomalies.

L.4. VLMs cannot count rows and columns in simple board game grids 🕮

Experiments Following our previous tasks, we evaluate five VLMs on counting tasks in four \boxplus grid-based game boards: chess (8×8), Go (19×19), Sudoku (9×9), and xiangqi (10×9). For chess (see Fig. F20e) and Sudoku (see Fig. F20c), models are asked to report the number of rows and columns. For Go and xiangqi (see Fig. 3f), they report the counts of horizontal and vertical lines.

Results All VLMs perform extremely poorly on \mathbb{H} board game grid counting, (2.26% mean accuracy; Tab. T4). The models even failed to answer any counting questions correctly on Sudoku (see Fig. F20c) and Go (0%; Tab. T4). These findings confirm that current VLMs are unable to execute basic visual counting tasks in structured settings and instead default to overconfident but incorrect guesses.

M. Prompts used for image generation and image editing

Table T5: Prompts used for image generation and image editing with +/Gemini-2.0 Flash and GPT-40 by topic and prompt type

Topic	Prompt type	Prompt
Animals	Animal suggestions	Generate a JSON list containing 100 animal objects. Each object should represent a common animal and follow the structure below: { "name": " <common animal="" name="">", "num_legs": <typical legs="" number="" of=""> } Ensure the following for each animal: 1. the number of legs of this animal is 2 or 4. 2. the animal's legs must be long enough to be seen easily from the body using a side-view perspective. Prioritize animals whose legs are thin and/or long.</typical></common>
	Animal generation	Generate a clear, full-body, side-view image of a(n) {animal} with {num_legs} legs that is walking in a real-world natural background. The {num_legs}-legged animal must look photo-realistic in nature. All {num_legs} legs must be clearly visible.
	Animal editing	Edit this image: Add 1 more leg to the {animal} so that it has {num_leg} legs in total. The {num_leg}-legged {animal} must be photo-realistic. All {num_leg} legs must be clearly visible.
Flags	Flag suggestions	Generate a JSON list of flags objects. Each object should represent a well-known flags and follow the structure below: { "name": " <flag name="">", "original_stripes" or "original_stars": <number (whichever="" applicable)="" of="" or="" stars="" stripes=""> } 1. Ensure that the number of stars is more than 3, and the number of stripes is at least 5. 2. Ensure that the flag does not contain any other geometrically complex elements (depicting of animal, letters, etc.). 3. Prioritize well-known flags.</number></flag>
	Flag SVG code editing	You are an expert in editing SVG image code. Modify the SVG code of the flag of {country} according to the following instruction: Instruction: "The flag of {country} has {num_ele} {element}. Modify the SVG code so that it has num_ele + 1 {element} instead. Make sure the modified {element} are natural looking and integrate seamlessly on the new flag." Base SVG code: {svg_code} 1. Modify the base SVG by adding or removing the mentioned feature (stars, stripes, etc.) according to the instruction above. 2. Wrap the entire SVG in <code></code> . Do not explain anything.

Table T6: Prompts used for image generation and image editing with +/Gemini-2.0 Flash and @GPT-40 by topic and prompt type

Topic	Prompt type	Prompt					
Logos	Logo suggestion	Generate a JSON list of subtle logo modification prompts and correspond- ing VLM question prompts to test visual bias. For each entry: Slightly modify the visual components of a well-known car or sportswear logo. The selected logo must be geometrically simple and widely recognized. You must include a generation prompt to create the altered image. Include a question prompt (e.g., "How many"). Include metadata: element being modified, actual count (ground truth), common expected count (bias). <in-context 1="" example="" learning=""> <in-context 2="" example="" learning=""></in-context></in-context>					
	Shoe generation	Generate an {shoe_brand} style running shoe but with {actual_count} {modified_element} instead of {expected_bias}.					
	Shoe background generation	Generate a side-view image of an athlete wearing this pair of shoes. Keep all the fine-grained details of the shoes, particularly the {actual_count} {modified_element} on both shoes. The per- son is playing {sports_type}, showing their sports_type skills, and is wearing a {sports_type} outfit. Zoom out a bit to see their full body.					
	Car logo generation	Generate a {car_brand} logo but with {actual_count} {modified_element} instead of {expected_bias}.					
	Car background generation	Generate a photo-realistic front-view image of a $\{color\}$ $\{car_brand\}$ $\{body_type\}$ on the road in the middle of the day. Zoom out a bit so that we can see the road.					

N. Questions for sanity check

Торіс	Identification questions	Counting/Illusion questions
Animal	What animal is this? Answer in curly brackets, e.g., {Fish}.	How many legs do this animal have? Answer with a number in curly brackets, e.g., {9}.
Logo	What car logo is this? Answer in curly brackets, e.g., {Toyota}.	How many overlapping circles are there on the logo of this car? Answer with a number in curly brackets, e.g., {9}.
Flags	What country flag is this? Answer in curly brackets, e.g., {Flag of Vietnam}.	How many stars are there in this flag? Answer with a number in curly brackets, e.g., {9}.
Chess Pieces	What board game is this? Answer in curly brackets, e.g., {Shogi}.	How many chess pieces are there on this board? Answer with a number in curly brackets, e.g., {9}.
Game Boards	What board game is this? Answer in curly brackets, e.g., {Shogi}.	How many rows are there on this board? Answer with a number in curly brackets, e.g., {9}.
Optical Illusions	What optical illusion is this? Answer in curly brackets, e.g., {Delboeuf illusion}.	This image shows the Ebbinghaus illusion. What ques- tion does this illusion typically ask, and what is the correct answer?

Table T7: Examples of Sanity check questions