# BLINK👁: Multimodal Large Language Models Can See but Not Perceive

Xingyu Fu[1][⋆], Yushi Hu[2,3][*], Bangzheng Li[4], Yu Feng[1], Haoyu Wang[1], Xudong Lin[5], Dan Roth[1], Noah A. Smith[2,3], Wei-Chiu Ma[3][†], Ranjay Krishna[2,3][†]

[1]University of Pennsylvania,[2]University of Washington,[3]Allen Institute for AI,
[4]University of California, Davis,[5]Columbia University
https://zeyofu.github.io/blink/

**Figure 1: The BLINK Benchmark.** BLINK contains 14 visual perception tasks that can be solved by humans "within a blink", but pose significant challenges for current multimodal LLMs. These tasks are inspired by classical computer vision problems and recast into multiple-choice questions for multimodal LLMs to answer. Notice that the visual prompts and questions in this figure are different from the actual ones used in the benchmark for illustrative purposes, and answers of the samples are provided.[1]

**Abstract.** We introduce BLINK, a new benchmark for multimodal language models (LLMs) that focuses on core visual perception abilities not found in other evaluations. Most of the BLINK tasks can be solved by humans "within a blink" (*e.g.*, relative depth estimation, visual correspondence, forensics detection, and multi-view reasoning). However, we find these perception-demanding tasks cast significant challenges for current multimodal LLMs because they resist mediation through natural language. BLINK reformats 14 classic computer vision tasks into 3,807 multiple-choice questions, paired with single or multiple images and visual prompting. While humans get 95.70% accuracy on average,

---

[⋆] Equal contribution. Correspond to <Xingyu Fu: xingyuf2@seas.upenn.edu>, <Yushi Hu: yushihu@uw.edu>. All data and evaluation are available on the project page.

[†] Both authors advised equally.

Blink is surprisingly challenging for existing multimodal LLMs: even the best-performing GPT-4V and Gemini achieve accuracies of 51.26% and 45.72%, only 13.17% and 7.63% higher than random guessing, indicating that such perception abilities have not "emerged" yet in recent multimodal LLMs. Our analysis also highlights that specialist CV models could solve these problems much better, suggesting potential pathways for future improvements. We believe Blink will stimulate the community to help multimodal LLMs catch up with human-level visual perception.

## 1   Introduction

Compared to today, computer vision was originally attempting to interpret images as projections of 3D scenes, not just processing 2D arrays of flat "patterns" [25, 58, 61]. In this pursuit, early research developed a series of intermediate tasks: they focused on understanding optical properties like reflectance [12, 76], 3D primitives through multi-view reasoning [37, 59], geometric reasoning through depth estimation [74], instance recognition through visual correspondence [55], affordance through keypoint grounding [36], and forensics through intrinsic images [9]. Yet in the modern era of large language models (LLMs), we, as a community, have focused less on such perceptual tasks, and instead have developed new tasks, mostly expressed in natural language, emphasizing the vision-language connection learned by multimodal LLMs [3, 6, 18, 20, 24, 26, 50, 51, 56, 62, 71, 77]. This might be because many traditional computer vision tasks resist mediation through natural language, due to the inherent imprecision of language (*e.g.*, it is challenging to precisely pinpoint a spatial keypoint through language).

This paper aims to highlight crucial aspects of visual perception that have been overlooked when evaluating multimodal LLMs. To appropriately position our paper, let us revisit how we currently evaluate perception through using multimodal LLMs [43, 44, 47, 52, 53, 57, 87]. While many of these benchmarks have been popularized as the de facto evaluation measures for influential models like GPT-4V and Gemini-Pro, they conflate perception with language knowledge and reasoning. At the risk of singling out one benchmark, let us consider two questions highlighted in the popular MMBench [52]: "`<image 1> Why is this hummingbird called ruby-throated?`" and "`<image 1> What will happen next? A: the person is gonna laugh B: the person is gonna cry.`" For the first question, the vision subpart is to recognize the hummingbird. For the second, it only needs a coarse description of the image. Everything else is left to the language model to solve. Such a conflation has also been reported for other benchmarks by previous work [11, 38, 84]. Our experiments show that this conflation reductively evaluates perception as a dense captioning task. In other words, **by replacing the image with a task-agnostic dense caption, our experiments show that a "blind" GPT-4 performs well on these "multimodal tasks".**

In response, we propose Blink. Blink reimagines traditional computer vision problems through a format that allows us to evaluate multimodal LLMs.

**Figure 2: Comparison between BLINK and previous benchmarks.** BLINK has several novel features: (1) BLINK incorporates diverse visual prompts, like circles, boxes, and image masks, while previous benchmarks only have text questions and answers. (2) BLINK evaluates a more comprehensive range of visual perception abilities, like multi-view reasoning, depth estimation, and reflectance estimation. Prior benchmarks are generally more focused on recognition-based VQA. (3) BLINK contains "visual" commonsense problems that humans can answer within seconds, while prior benchmarks like [87] require domain knowledge. The samples of previous benchmarks are from [44, 52, 87]. Part of our samples are curated from [10, 19, 29, 31, 35, 42, 88].

As partially demonstrated in Figure 1,[1] BLINK consists of 14 classic computer vision tasks, ranging from low-level pattern matching (*e.g.*, visual correspondences estimation) to mid-level reasoning (*e.g.*, relative depth estimation), and extending to high-level visual understanding (*e.g.*, visual similarity). The image tasks are meticulously selected such that they are difficult to solve by reducing the evaluation using dense captioning; instead, the models must perceive the contents of the image(s) to answer. We recast each traditional task into a modern question-answering format, where answer choices are either images or text. BLINK contains 3.8K questions across 7.3K images, where questions may contain multiple images that are curated from a wide range of datasets [8, 10, 19, 29, 31, 35, 41, 46], encompassing indoor household scenes as well as outdoor urban or natural environments. The questions and choices are either derived from the datasets, or

---

[1] The answers of the examples in Figure 1 are as follows. Relative depth: B; jigsaw: A; multi-view reasoning: right; visual correspondence: A; semantic correspondence: C; forensics detection: final image; IQ test: D; visual similarity: upper one; functional correspondence: A; relative reflectance: they are about the same.

manually written by humans. On average, each question can be solved by a human subject within a BLINK of an eye, except the IQ test.

We carefully evaluate 17 multimodal LLMs with various sizes (*i.e.*, 7B, 13B, 34B) on BLINK. We observe the paradox that **while these problems are easy for humans (**95.70% **average accuracy), they are extremely hard for existing machinery** – even GPT-4V model can only achieve 51.26% accuracy on average, which is 44.44% worse than humans, and 13.17% better than random guessing. We also experiment with specialist vision models and find that they perform much better than multimodal LLMs. For example, the specialist outperforms GPT-4V by 62.8% on visual correspondence estimation, 38.7% on relative depth estimation, and 34.6% on multi-view reasoning, in terms of absolute accuracy. Our findings indicate that the perceptual abilities of multimodal LLMs have been previously overestimated. Furthermore, these models may benefit from integrating insights from specialized models that excel in these areas. We believe BLINK can serve as an effective testbed for bridging the gap between traditional notions of perception and the modern generative capabilities of multimodal LLMs.

## 2   Related Work

### 2.1   Multimodal Models

Inspired by the impressive success in recent large language models (LLMs) [13, 21, 62, 75, 90], a sequence of studies explore multimodal LMMs that can jointly understand vision and language information and generate textual answers through adding a modality adaption structure between a frozen visual encoder [27, 64, 69] and a frozen LLM [75, 90]. Flamingo [3] and BLIP-2 [45] are two of the earliest works to explore these transformer-based multi-modality conjunction structures. They first pre-train on image-text matching datasets [15, 41, 46, 65] and then fine-tune on task-specific datasets such as visual question answer (VQA) [4, 33]. Starting from LLaVA [48, 50], people use LLM synthesized instruction-following chat data (which are in VQA format) for instruction tuning and achieve much better results [7, 16, 24, 71]. There have been extended studies that explore further capabilities of multimodal LLMs, especially on VQA reasoning [30–32, 38, 39, 66, 81, 89]. However, they mainly focus on the textual reasoning abilities [80] within the multimodal LLMs and do not emphasize visual perceptions.

### 2.2   Multimodal Benchmarks

Traditional vision-language datasets are designed to assess single-task capabilities, such as optical character recognition (OCR) [54], image captioning [46], and visual question answering [4, 33]. However, these datasets are often not comprehensive enough to holistically assess multimodal LMMs on general perception and reasoning abilities. Many recent papers have built more comprehensive benchmarks. MME [28] is one of the earliest holistic benchmarks containing multi-modal Yes/No questions on the defined visual perception and language

reasoning tasks. MM-Vet [86] includes six sub-features from the previous datasets including recognition-focused questions, OCR, and math, providing a diverse while discrete evaluation set. MMBench [52] covers more subjects and provides a more robust circular evaluation setting. Seed-Bench [43, 44] benchmark has a more diverse source of inputs, including multiple-image inputs and video, and includes more tasks. However, the visual perception questions in MME, MMBench, MM-Vet, and Seed-Bench are mainly extracted from existing VQA datasets or generated by GPT [62] from image descriptions such as COCO-Caption [46], and are recognition focused, covering topics such as object (attribute)recognition, and OCR. In contrast, we focus on multiple distinct nuanced perception abilities and recognition-level perception is only one of our focus. Some other multimodal benchmarks have distinct focuses. MMMU [87] aims at achieving expert-level artificial general intelligence by collecting domain-knowledge-required questions. HallusionBench [34] mainly tests the language hallucination and visual illusion phenomena. MathVista [57] presents exclusively mathematical domain visual questions based on images such as charts, tables, and diagrams. These benchmarks do not require human-level perception abilities as in BLINK and therefore cannot measure model visual perceptions holistically.

## 3   The BLINK Benchmark

Our goal is to faithfully evaluate the visual perception capabilities of existing Multimodal LLMs. We seek to study the visual perception gap between humans and machineries, and offer deeper insights into potential pathways towards achieving more generalized machine perception. Based on the observation that existing benchmarks predominantly focus on evaluating visual recognition abilities, we introduce a novel benchmark, BLINK, designed to enable both quantitative and qualitative evaluation of the nuanced perception capabilities of multimodal LLM across various dimensions. We unfold this section by illustrating the overall design of BLINK (§3.1) and discussing its unique features comparing with previous benchmarks. Then we describe each task in detail, providing an in-depth explanation of the data curation process (§3.2).

### 3.1   Overview of BLINK

To ensure that one can effectively measure what Multimodal LLMs can or cannot perceive, we carefully select 14 tasks (see §3.2 for the full list) that are difficult to solve by reducing the evaluation into text-only questions using dense captioning. The tasks are drawn from either classic computer vision problems or recent applications of Multimodal LLMs, each of which requires a nuanced understanding of the visual data. They range from low-level pattern matching (e.g., visual correspondence) to mid-level spatial reasoning (e.g., relative depth), and up to high-level visual understanding (e.g., visual similarity). This variety allows for a systematic exploration of Multimodal LLMs' capabilities across different perceptual complexity layers. Furthermore, these visual tasks vary in
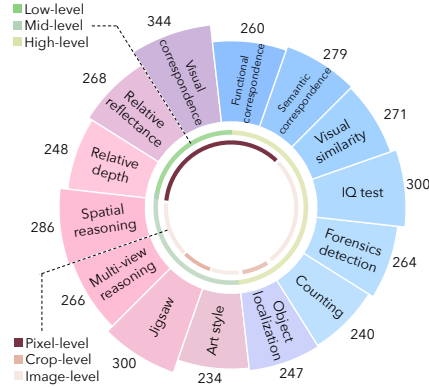
**Figure 3:** Statistics of BLINK. The benchmark includes 14 tasks, ranging from pixel-level to image-level perception, and low-level pattern matching (*e.g.*, visual correspondences estimation) to mid-level reasoning (*e.g.*, relative depth estimation), and extending to high-level visual understanding (*e.g.*, visual similarity).
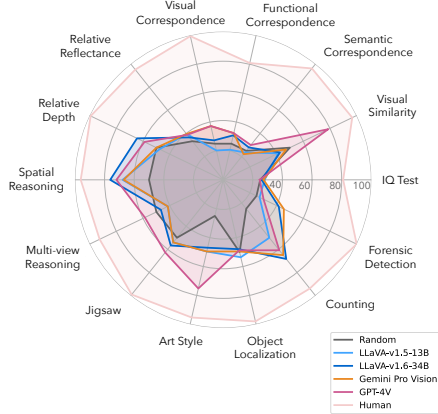


**Figure 4:** Accuracies of multimodal LLMs on BLINK test set. Please refer to Table 1 and §4.2 for more results and discussions.

granularity, ranging from pixels (*e.g.*, relative reflectance) to patches (*e.g.*, jigsaw) and extending to the full image (*e.g.*, forensic detection), enabling us to evaluate models' proficiency in observing at various scales.

To facilitate the evaluation of multimodal LLMs, we recast all tasks as multiple-choice question-answering problems. The options for answers may include images or texts, while the questions themselves can feature either single or multiple images. Prompts are designed to be both textual and visual in nature. We re-purposed several existing vision datasets as well as collected new data. In total, we contribute 3.9K multiple-choice questions and 7.3K images, with an even distribution between the validation and test sets. Numbers of each task are reported in Figure 3, and more detailed statistics can be found in Appendix A.5.

**Key features of BLINK:** Comparing with previous benchmarks, BLINK has the following novel features:

- **Visual prompting**: Unlike existing benchmarks that support only text prompting, BLINK features a variety of visual prompts. This enables one to highlight specific areas within images, facilitating the evaluation of Multimodal LLMs' detailed understanding of these regions. It also offers an interface for researchers to investigate the impact of visual prompting techniques.
- **Perception beyond recognition**: Besides visual recognition, BLINK considers a diverse set of visual perception abilities, such as 3D reasoning, geometric understanding, affordance reasoning, etc. The breadth allows one to evaluate Multimodal LLMs from an unique array of perspectives.

– **"Visual commonsense" that does not require domain knowledge**: The questions in BLINK are intentionally designed to be straightforward, requiring neither domain-specific knowledge nor expertise to answer. They are crafted in such a way that humans can solve them almost instantaneously, typically within a few seconds. This allows us to explore the fundamental gap in visual perception gap between humans and Multimodal LLMs, highlighting the paradox that problems easily solved by humans often pose significant challenges for machines.
– **Interleaved image-text formats**: BLINK features a heterogeneous question-answering format, wherein both questions and choices can be presented as text or images. This diversity compels Multimodal LLMs to genuinely understand the questions, pushing the boundaries of their interpretative capabilities.
– **Diverse image sources**: BLINK comprises a wide range of in-the-wild images sourced from various origins, covering everything from indoor and outdoor scenes to object-centric views and landscapes. This collection spans abstract diagrams, synthesized images, and authentic photographs, ensuring a comprehensive examination of visual perception

The design principles of BLINK are also illustrated in Figure 2. We will now describe each task in detail.

## 3.2 Dataset Collection Process

BLINK comprises 14 tasks, all of which have been repurposed into a multiple-choice question-answering format. These tasks utilize a diverse collection of images from various sources, and we ensure that each test sample across all tasks features unique images.

**Visual correspondence:** This task aims to evaluate the ability of Multimodal LLMs to understand and identify the same scene point across various viewpoints, lighting conditions, or time. We exploit HPatches [8] for this task. HPatches contains a number of image sequences, each of which are composed of images taken under different illuminations and/or viewpoints of a scene. For each question, we randomly sample two images and an interest point within them. Then we exploit the ground-truth homography to compute its correspondence. Finally, we randomly select three more interest points to serves as other choices.

**Relative reflectance:** This task aims to compare the reflectance (albedo) of two pixels. It allows us to evaluate Multimodal LLMs' understanding of material properties and their interaction with light, which is crucial for applications requiring high-fidelity visual interpretations. We curate our samples using human annotations from the Intrinsic Images in the Wild (IIW) dataset [10]. Each question is based on an image and two specified points, with the objective being to identify which point is darker, or whether the two points have similar reflectance.

**Relative depth:** Humans are good at judging relative depth [19]. This task can thus serve as a proxy to validate if the geometric understanding capabilities

of existing multimodal LLMs are close to human. We curate our samples using human annotations from the Depth in the Wild [19] dataset. Each question contains an image and two specified points. The task is to determine which point is closer.



**Figure 5: Qualitative results** on BLINK. For each task, we show the choice of LLaVA-v1.6-34B [50], Qwen-VL-Max [7], Gemini Pro [71], GPT-4V [62], and humans. Red choice indicates the ground truth. Notice that the markers are intentionally enlarged for visualization purposes, and we make some images inset images to save space. For IQ test, the third image is constructed by overlaying the first and second images.

**Spatial relation:**  Understanding spatial relationships between objects in a scene is essential for interpreting complex visual environments. However, modern Multimodal LLMs often struggles with spatial concepts such as "left" and "right" [85]. This task help us evaluate whether the models finally possess this vital skill. We curate our samples from the Visual Spatial Reasoning [47] dataset. Each sample contains an image and a claim. The task is to determine if the claim is true or false. We reformat the claims into binary questions via GPT-3.5 [13].

**Multi-view reasoning:** This task is centered on evaluating the multi-view reasoning capabilities of Multimodal LLMs. The objective is to deduce the relative camera motion based on two images of an object captured from different viewpoints. Our data is sourced from the Wild6D dataset [88], which features videos of various objects recorded in diverse settings. We select two random frames from each video to calculate the relative camera motion. Recognizing that even humans might struggle to precisely articulate 3D motion details, we simplify the task by classifying motions into two broad categories: moving towards the left or moving towards the right. Despite the simplicity of these questions, as we will later demonstrate, they pose significant challenges for current models.

**Jigsaw:**  This task assesses the ability of Multimodal LLMs to recognize and group patterns, as well as to align patches based on continuity in shape, color, and texture. We utilize images from the TARA dataset [31] and segment each of them into a 3x3 grid. We retain the three segments from the upper left corner as the reference image, and treat the central segment along with a randomly chosen segment as options. The objective is to identify the correct patch (*i.e.*, the central patch).

**Art style:**  This task evaluates Multimodal LLMs capability to analyze and discern both local and global similarities in art styles among multiple images. Although there have been prior efforts to incorporate art-related questions into evaluation [87], such attempts primarily focused on questions requiring expert-level knowledge, including deducing an artist's name and understanding historical contexts, rather than on direct image comparison. For this task, we collect paintings and their stylistic information from WikiArt. Given one reference painting image and two other paintings as options, the model is tasked with identifying the one that most closely shares the art style of the reference painting.

**Object localization:**  The ability to accurately detect and localize objects is critical for scene understanding. While previous benchmarks [52] have explored this task, their focus was primarily on coarse localization. For instance, they might only ask the model if an object is located at the "top" or "right" side of an image. BLINK, in contrast, aims for a more fine-grained evaluation. We exploit images from LVIS [35], randomly sampling one object per image along with its ground-truth bounding box. Then we add Gaussian noise to the ground-truth box to create a confounding box. The goal is to select the correct one.

**Counting:**  This task evaluate Multimodal LLMs' abilities in detection, recognition, and compositional reasoning, particularly in complex scenes where objects may overlap, be occluded, or vary in size and appearance. We select our questions

from the TallyQA dataset [2], known for its challenging human-written counting questions. Each sample comprises an image, a question, and a numerical answer. In addition to the correct answer, we randomly select three numbers to serve as confounding options.

**Forensic detection:** Recent advances in generative AI have raised concerns about malicious uses and have prompted calls for the automatic detection of fake content. To evaluate whether Multimodal LLMs can fulfill such a role, we construct sets of real and synthesized images that describe similar scenes and ask the models to identify the real ones. Specifically, we first generate synthetic images using Stable Diffusion XL [63], employing COCO captions [46] as prompts. Then, we manually search online using these captions as descriptions and select high-quality photographs as the real images.

**IQ test:** This task evaluates the ability of Multimodal LLMs to engage in graphical reasoning, without requiring any domain-specific knowledge. We manually collect test samples, along with human explanations, from various public, license-friendly online sources. Given visual examples and a selection of images, the objective is to identify the image that either continues the pattern established by the examples or is spatially consistent with them.

**Visual similarity:** This task aims to verify whether Multimodal LLMs possess a nuanced understanding of visual features, patterns, and aesthetics at a level comparable to humans. We select our samples from the DreamSim dataset [29]. Given a reference image alongside two alternative images, the objective is to identify the image that most closely resembles the reference image in terms of visual similarity.

**Semantic correspondence:** This task focuses on identifying and matching semantically similar yet visually distinct elements across images, thereby evaluating the ability of Multimodal LLMs to understand the underlying semantics of object parts. Our samples are sourced from the SPair-71k dataset [60], which features pairs of images with multiple corresponding semantic points. For each task, we randomly select one semantic point in an image as a reference, and provide the matching point alongside three random semantic points in the paired image as options. The objective is to accurately identify the correct matches.

**Functional correspondence:** The task aims to identify points that are functionally similar across objects. It challenges Multimodal LLMs to extend their understanding beyond mere semantics, enabling them to infer the diverse functions an object can perform in various contexts. Such capability is crucial for applications in robotics. We derive our samples from the FunKPoint dataset [42], which features paired images annotated for functional correspondences. Following a method analogous to semantic correspondence, we present an action alongside two object images. One image includes a reference point, while the other offers four potential points. The objective is to select the point that best matches the reference in terms of functional affordances.

**Data quality control:** To guarantee the quality of Blink, we manually go through all collected data and filter out data that are ambiguous.

## 4  Experiments

In this section, we first describe the experimental setup and the baselines (§4.1). Then we present a comprehensive evaluation of 16 recent multimodal LLMs (§4.2). We demonstrate that while humans can answer the questions with high accuracy, Blink is challenging for existing models. Finally, we provide detailed analyses on multiple experimental settings, including the effect of reducing images to captions, sensitivity to different visual prompts, and error analysis (§4.3).

### 4.1  Experimental Setup

**Multimodal LLMs:** We evaluate Blink on 16 recent multimodal LLMs, including MiniGPT-4-v2 [16], OpenFlamingo-v2 [5], InstructBLIP (7B and 13B) [24], CogVLM [77], LLaVA(v1, v1.5, v1.6, internLM, and xtuner versions, model size 7B, 13B, and 34B) [23, 26, 49–51], Yi-VL (6B and 34B)[2], Qwen-VL-MAX [7], Gemini Pro [71], Claude 3 Opus [1] and GPT-4V(vision) [62]. See Appendix B for more details.

**Evaluation setup:** We follow standard setups as in the VLMEvalKit [22], where the temperature is set to 0 and retry is set to 10. However, we do not resize the images during any experiment. For the models that do not support multiple images as input, we concatenate the images as input. We extract the choice from the models' output with a set of pre-defined rules and GPT-3.5-turbo [13]. We refer the readers to Appendix A for more details on visual prompting, how we generate the answers in Blink, and the human evaluation protocol.

### 4.2  Main Results

**Overall performance:** As shown in Table 1, the mean accuracy of 7B and 13B open-source Multimodal LLMs hover around 35–42%, which is similar to random guess (38.09%). The most proficient open-source model, LLaVA-v1.6-34B, achieves an accuracy of 45.05%. Even the most advanced models, GPT-4V and Gemini Pro and Claude 3 OPUS, achieve accuracies of only 51.26%, 45.72%, and 44.11% respectively. Their performance are merely 13.17%, 7.63% and 6.02% better than random guessing and lag behind human performance by 44.44%, 49.98% and 51.59%. Notably, for certain tasks such as jigsaw, semantic correspondence, multi-view reasoning, object localization, and relative reflectance, some multimodal LLMs even underperform compared to random guessing. Some qualitative results are shown in Figure 5.

**In which tasks do multimodal LLMs show relative strengths and weaknesses?** Figure 4 shows the accuracies of the best-performing models on Blink: LLaVA-v1.6-34B [50], Gemini Pro [71], and GPT-4V [62]. We observe that multimodal LLMs perform relatively better on spatial reasoning, art style, and counting tasks, in which they are much better than random guessing. The

---

[2] More details are at the official website at https://www.01.ai/

| | Validation (1,901) | Test (1,906) | Similarity (136) | Counting (120) | Depth (124) | Jigsaw (150) | Art (117) | Fun.Corr. (130) |
|---|---|---|---|---|---|---|---|---|
| Random Choice | 38.09 | 38.09 | 50 | 25 | 50 | 50 | 50 | 25 |
| Human | 95.67 | 95.70 | 96.70 | 93.75 | 99.19 | 99.00 | 95.30 | 80.77 |
| **Open-source multimodal LLMs** | | | | | | | | |
| MiniGPT-4-v2 [16] | 34.23 | 34.57 | 52.94 | 10.83 | 49.19 | 26.00 | 47.86 | 18.46 |
| OpenFlamingo-v2 [5] | 39.18 | 38.32 | 55.15 | 21.67 | 54.03 | 46.00 | 52.14 | 36.15 |
| InstructBLIP-7B [24] | 39.72 | 38.65 | 46.32 | 29.17 | 50.81 | 54.00 | 47.86 | 23.85 |
| InstructBLIP-13B [24] | 42.24 | 39.58 | 46.32 | 30.83 | 50.00 | 54.00 | 50.43 | 22.31 |
| LLaVA-internLM2-7B [72] | 37.71 | 36.06 | 52.94 | 52.50 | 52.42 | 34.67 | 30.77 | 23.08 |
| Yi-VL-6B [2] | 38.72 | 41.24 | 46.67 | 46.67 | 56.45 | 50.00 | 53.85 | 23.85 |
| Yi-VL-34B [2] | 41.68 | 42.78 | 50.00 | 58.33 | 53.23 | 54.00 | 46.15 | **39.23** |
| LLaVA-v1.5-7B-xtuner [23] | 39.36 | 40.81 | 46.32 | 53.33 | 50.81 | 54.00 | 47.86 | 23.85 |
| LLaVA-v1.5-13B-xtuner [23] | 42.00 | 41.31 | 46.32 | 45.00 | 54.03 | 53.33 | 47.86 | 26.15 |
| CogVLM [77] | 41.54 | 39.38 | 46.32 | 38.33 | 50.81 | 52.67 | 49.57 | 23.85 |
| LLaVA-v1.5-7B [48] | 37.13 | 38.01 | 46.32 | 43.33 | 50.81 | 11.33 | 47.86 | 21.54 |
| LLaVA-v1.5-13B [48] | 42.66 | 40.55 | 46.32 | 50.00 | 47.58 | 54.00 | 47.86 | 20.77 |
| LLaVA-v1.6-34B [50] | 46.80 | 45.05 | 46.32 | **68.33** | 64.52 | 56.67 | 47.01 | 30.77 |
| **API-based models** | | | | | | | | |
| Qwen-VL-Max [7] | 40.28 | 41.94 | 51.47 | 55.83 | 58.87 | 3.33 | 37.61 | 28.46 |
| Gemini Pro [71] | 45.16 | 45.72 | 55.88 | 65.00 | 50.00 | 54.00 | 49.57 | 32.31 |
| Claude 3 OPUS [1] | 44.05 | 44.11 | 70.59 | 49.17 | 57.26 | 32.67 | 60.68 | 22.31 |
| GPT-4V(ision) [62] | 51.14 | 51.26 | **83.09** | 60.83 | 58.87 | 62.67 | 78.63 | 31.54 |
| GPT-4 Turbo [62] | 54.61 | 53.89 | **83.09** | 60.83 | **66.94** | **66.00** | 81.20 | 31.54 |
| GPT-4o [62] | **60.04** | **59.03** | 65.44 | 51.67 | 64.52 | 58.00 | **82.91** | **39.23** |

| | Sem.Corr. (140) | Spatial (143) | Local. (125) | Vis.Corr. (172) | Multi-view (133) | Reflect. (134) | Forensic (132) | IQ (150) |
|---|---|---|---|---|---|---|---|---|
| Random Choice | 25 | 50 | 50 | 25 | 50 | 33.33 | 25 | 25 |
| Human | 96.07 | 98.25 | 98.00 | 99.42 | 92.48 | 95.14 | 100.00 | 80.00 |
| **Open-source multimodal LLMs** | | | | | | | | |
| MiniGPT-4-v2 [16] | 26.43 | 51.75 | **56.00** | 23.84 | 52.63 | 31.34 | 17.42 | 19.33 |
| OpenFlamingo-v2 [5] | 23.57 | 46.85 | 52.00 | 25.00 | 41.35 | 43.28 | 15.91 | 23.33 |
| InstructBLIP-7B [24] | 25.00 | 55.24 | 44.80 | 22.67 | **58.65** | 29.85 | 29.55 | 23.33 |
| InstructBLIP-13B [24] | 22.86 | 64.34 | 52.00 | 20.93 | 54.14 | 46.27 | 13.64 | 26.00 |
| LLaVA-internLM2-7B [72] | 22.14 | 74.13 | 48.00 | 21.51 | 41.35 | 32.84 | 3.79 | 14.67 |
| Yi-VL-6B [2] | 26.43 | 72.73 | 49.60 | 29.65 | 48.12 | 29.85 | 20.45 | 23.33 |
| Yi-VL-34B [2] | 21.43 | 70.63 | 54.40 | 23.84 | 41.35 | 46.27 | 17.42 | 22.67 |
| LLaVA-v1.5-7B-xtuner [23] | 24.29 | 74.83 | 45.60 | 23.84 | 42.11 | 26.87 | 36.36 | 21.33 |
| LLaVA-v1.5-13B-xtuner [23] | 22.14 | **77.62** | 48.00 | 22.09 | 41.35 | 46.27 | 29.55 | 18.67 |
| CogVLM [77] | 23.57 | 67.13 | 43.20 | 20.93 | 57.14 | 26.87 | 24.24 | 26.67 |
| LLaVA-v1.5-7B [48] | 32.14 | 70.63 | 48.80 | 20.35 | 49.62 | 36.57 | 28.03 | 24.00 |
| LLaVA-v1.5-13B [48] | 23.57 | 67.83 | 47.20 | 20.35 | 41.35 | 45.52 | 27.27 | 28.00 |
| LLaVA-v1.6-34B [50] | 27.86 | 76.22 | 41.60 | 27.33 | 46.62 | 29.85 | 41.67 | 26.00 |
| **API-based models** | | | | | | | | |
| Qwen-VL-Max [7] | 29.29 | **77.62** | 49.60 | 22.67 | 53.38 | **49.25** | 47.73 | 22.00 |
| Gemini Pro [71] | 22.14 | 67.13 | 46.40 | 37.21 | 41.35 | 46.27 | 45.45 | 27.33 |
| Claude 3 OPUS [1] | 20.71 | 57.34 | 46.40 | 31.40 | 57.89 | 27.61 | 62.12 | 21.33 |
| GPT-4V(ision) [62] | 30.00 | 72.03 | 50.40 | 37.21 | 58.65 | 38.81 | 30.30 | 24.67 |
| GPT-4 Turbo [62] | 32.86 | 67.13 | 48.80 | 42.44 | 57.14 | 34.33 | 51.52 | **30.67** |
| GPT-4o [62] | **45.71** | 76.92 | **56.00** | **71.51** | **60.15** | 38.81 | **85.61** | 30.00 |

**Table 1: Results of different models on the BLINK test set**. The first row shows task names and number of test data. The best performance in each task is in-bold. For the sake of completion, we also show the average score on the validation set. Detailed scores on the validation set are in Appendix C.

models also demonstrate some capability in relative depth and forensics detection. Overall, they are doing relatively well on mid-level perception tasks. In terms
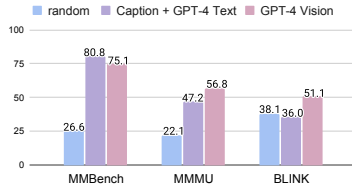
**Figure 6:** Performance of using image caption + text-only GPT-4 *vs.* GPT-4 Vision on MMBench [52], MMMU [87], and BLINK (§4.3).
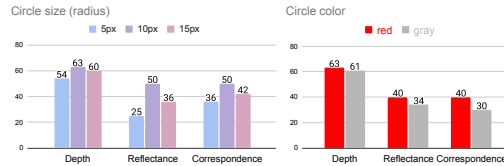


**Figure 7:** Accuracy of GPT-4V with different visual prompts (*e.g.*, different circle sizes, colors) on relative depth, relative reflectance, and visual correspondence tasks. More discussions in §4.3.

of granularity, the models in general perform better on image-level tasks and struggle on pixel-level and crop-level tasks.

**GPT-4V behaves differently:** Figure 4 and Table 1 show an interesting phenomenon: GPT-4V's performance pattern is different from other models. Compared with its counterparts, GPT-4V is much better in visual similarity, art style, jigsaw, and multi-view reasoning. Specifically, its performance on visual similarity is 29% better than Gemini Pro, demonstrating that GPT-4V possesses a nuanced understanding of visual patterns and aesthetics that is similar to humans. In contrast, Gemini Pro and LLaVA have similar performance patterns.

**Human performance:** Human evaluators achieve over 95% accuracy across most tasks, with an average accuracy of 95.70% .[3] This performance disparity between humans and multimodal LLMs highlights the significant visual perception gap that exists between current machine learning models and humans in perceiving, processing, and understanding complex visual and textual context.

### 4.3   Analysis

**Is dense captioning all you need for a multimodal LLM benchmark?** To answer the question, we reduce multimodal benchmarks to a text-only problem. Specifically, we convert images into task-agnostic dense image captions with GPT-4V. The dense caption describes detailed information about the image and the visual prompts (*e.g.*, where each circle is), using language. For each multimodal question, we prompt the text-only GPT-4-0125-preview model with image captions and the textual question and evaluate if the "blind" GPT-4 can answer the question. We call this `Caption + LLM`. This experiment is predicated on the hypothesis that captioning involves predominantly recognition-centric perception. If using captions along with text-only LLMs yields performance comparable to or surpassing that achieved through the integration of images with multimodal LLMs, then the perception demands of that benchmark are primarily confined to recognition only.

We experiment with BLINK, MMBench [52] and MMMU [87], as illustrated in Figure 6. Surprisingly, we find that the `Caption + LLM` setting achieves better

---

[3] Note that the human score for IQ test is annotated by authors. It may not reflect typical human performance, which is also expected to vary.

| Task | Vis.Corr. | Depth | Multi-view | Sem.Corr. | Forensic | Reflect. |
|------|-----------|-------|------------|-----------|----------|----------|
| Random | 25.00 | 50.00 | 50.00 | 25.00 | 25.00 | 33.33 |
| Human | 99.56 | 99.59 | 92.10 | 94.60 | 100.00 | 99.63 |
| Gemini Pro | 42.44 | 40.32 | 44.36 | 26.62 | 50.76 | 45.52 |
| GPT-4V | 33.72 | 59.68 | 55.64 | 28.78 | 34.09 | 38.81 |
| Specialist | DIFT [70] | DepthAnything [83] | LoFTR [68] | DIFT [70] | DIRE [79] | Ordinal Shading [14] |
|  | 96.51 | 97.58 | 90.22 | 71.22 | 68.94 | 77.61 |

**Table 2:** Comparison between multimodal LLMs, specialists, and human performance on the BLINK dev set. The specialists perform much better than multimodal LLMs.

results on MMBench than GPT-4V (with 5.7% increase in accuracy). On MMMU, `Caption + LLM` achieves 47.2% accuracy, which is 9.6% lower than GPT-4V performance, but is still much better than random guessing. On BLINK, `Caption + LLM` fails, achieving random guessing performance. These results indicate that dense captions cover the visual information needed for MMBench. For MMMU, image captions carry a large portion of visual information needed to answer the domain-knowledge-specific questions. Meanwhile, the performance decrease observed in BLINK suggests the necessity for advanced perceptual abilities beyond what is currently attainable with general captions. This variance highlights the limitations of existing multimodal LLM benchmarks in addressing the full spectrum of visual perception.

**Effect of visual prompting on BLINK:** Several BLINK tasks involve visual prompting. Prior work [67] shows that factors like shape, size, and color may affect task performance, and circles give the best overall performance. Following [67], we adopt circles in BLINK and analyze the effect of circle sizes and colors on multiple tasks in Figure 7. We experiment with relative depth, relative reflectance, and visual correspondence, with 100 validation set samples per task. The images are all reshaped to 1024px height. We experiment with circles with 5px, 10px, and 15px radius, and with red or gray color. We find that red is better than gray for all tasks. Also, the optimal circle size is task-dependent. On average 10px circles work the best, and we use it for all evaluations in this paper. The experiments suggest that visual prompting can have a big impact on multimodal LLM performance, and improving visual prompts or improving model robustness to prompt variation is a promising direction for future research [82].

**Can specialist models solve BLINK tasks?** Specialists can serve as a proxy upper bound of how good multimodal LLMs could be. We download the trained checkpoints for six specialist models and evaluate them on BLINK. As shown in Table 2, the specialists perform much better than GPT-4V and Gemini Pro, outperforming the best multimodal LLM by 18% to 57% on these tasks. Specifically, DepthAnything [83] and DIFT [70] achieve human-level performance on depth estimation and visual correspondence, whereas multimodal LLMs fail miserably. This sheds light on the possibility that multimodal LLMs may progress on these tasks given the correct data and training strategy. For instance, one possible way is to distill existing specialist models into multimodal LLMs [40].

**Error analysis of GPT-4V:** We randomly sampled 140 error instances made by GPT-4V on BLINK, 10 per task, and meticulously examined them. The most common types of errors are: **Hallucinate fine-grained patterns and attributes** (24.2%): the model hallucinates the nuanced details of objects. This error is most common for relative reflectance, forensics detection, and jigsaw tasks. **Hallucinate visual prompt locations** (20.0%): the circle location described by the model is wrong. This is common for visual correspondence and relative depth tasks. Other errors include Failures on capturing overall setting or style (8.6%), and Failures on grounding an object (5.7%). More details are in Appendix C.3.

## 5    Conclusion

We introduced BLINK, a new multimodal LLM benchmark that evaluates core visual perception abilities not found in existing evaluations. While these tasks seem trivial for humans to solve "within a blink", we find they pose significant challenges for current multimodal LLMs. Even the powerful GPT-4V and Gemini models only achieve around 50% accuracy on Blink, far below the 95.7% human performance. We conduct extensive analysis, measuring the effect of converting images to dense captions, visual prompting, self-consistency, analyzing the capabilities of specialist models, and conducting error analysis. We highlight that specialist computer vision models are performing much better than GPT-4V and Gemini on BLINK, shedding light on the possibility that multimodal LLMs may have big progress on these tasks. Ultimately, Blink provides a simple yet effective testbed for multimodal LLMs to catch up with human-level visual perception.

## References

1. Introducing the next generation of claude. https://www.anthropic.com/news/claude-3-family (March 2024) 11, 12, 23, 24
2. Acharya, M., Kafle, K., Kanan, C.: Tallyqa: Answering complex counting questions. In: AAAI (2019) 10
3. Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al.: Flamingo: a visual language model for few-shot learning. Advances in Neural Information Processing Systems **35**, 23716–23736 (2022) 2, 4, 22
4. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: Vqa: Visual question answering. In: Proceedings of the IEEE international conference on computer vision. pp. 2425–2433 (2015) 4
5. Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa, S., Jitsev, J., Kornblith, S., Koh, P.W., Ilharco, G., Wortsman, M., Schmidt, L.: Openflamingo: An open-source framework for training large autoregressive vision-language models. arXiv preprint arXiv:2308.01390 (2023) 11, 12, 22, 24
6. Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., Zhou, J.: Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond (2023) 2
7. Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., Zhou, J.: Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. arXiv preprint arXiv:2308.12966 (2023) 4, 8, 11, 12, 23, 24

8. Balntas, V., Lenc, K., Vedaldi, A., Mikolajczyk, K.: Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors. In: CVPR (2017) 3, 7

9. Barrow, H., Tenenbaum, J., Hanson, A., Riseman, E.: Recovering intrinsic scene characteristics. Comput. vis. syst **2**(3-26),  2 (1978) 2

10. Bell, S., Bala, K., Snavely, N.: Intrinsic images in the wild. ACM Trans. on Graphics (SIGGRAPH) **33**(4) (2014) 3, 7

11. Berrios, W., Mittal, G., Thrush, T., Kiela, D., Singh, A.: Towards language models that can see: Computer vision through the lens of natural language. arXiv preprint arXiv:2306.16410 (2023) 2

12. Black, M.J., Anandan, P.: A framework for the robust estimation of optical flow. In: 1993 (4th) International Conference on Computer Vision. pp. 231–236. IEEE (1993) 2

13. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems **33**, 1877–1901 (2020) 4, 9, 11, 21, 22

14. Careaga, C., Aksoy, Y.: Intrinsic image decomposition via ordinal shading. ACM Trans. Graph. (2023) 14

15. Changpinyo, S., Sharma, P., Ding, N., Soricut, R.: Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In: CVPR (2021) 4

16. Chen, J., Zhu, D., Shen, X., Li, X., Liu, Z., Zhang, P., Krishnamoorthi, R., Chandra, V., Xiong, Y., Elhoseiny, M.: Minigpt-v2: large language model as a unified interface for vision-language multi-task learning (2023) 4, 11, 12, 24

17. Chen, J., Zhu, D., Shen, X., Li, X., Liu, Z., Zhang, P., Krishnamoorthi, R., Chandra, V., Xiong, Y., Elhoseiny, M.: Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. arXiv preprint arXiv:2310.09478 (2023) 22

18. Chen, L., Li, J., Dong, X., Zhang, P., He, C., Wang, J., Zhao, F., Lin, D.: Sharegpt4v: Improving large multi-modal models with better captions. arXiv preprint arXiv:2311.12793 (2023) 2

19. Chen, W., Fu, Z., Yang, D., Deng, J.: Single-image depth perception in the wild. Advances in neural information processing systems **29** (2016) 3, 7, 8

20. Chen, X., Djolonga, J., Padlewski, P., Mustafa, B., Changpinyo, S., Wu, J., Ruiz, C.R., Goodman, S., Wang, X., Tay, Y., Shakeri, S., Dehghani, M., Salz, D., Lucic, M., Tschannen, M., Nagrani, A., Hu, H., Joshi, M., Pang, B., Montgomery, C., Pietrzyk, P., Ritter, M., Piergiovanni, A., Minderer, M., Pavetic, F., Waters, A., Li, G., Alabdulmohsin, I., Beyer, L., Amelot, J., Lee, K., Steiner, A.P., Li, Y., Keysers, D., Arnab, A., Xu, Y., Rong, K., Kolesnikov, A., Seyedhosseini, M., Angelova, A., Zhai, X., Houlsby, N., Soricut, R.: Pali-x: On scaling up a multilingual vision and language model (2023) 2

21. Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H.W., Sutton, C., Gehrmann, S., Schuh, P., Shi, K., Tsvyashchenko, S., Maynez, J., Rao, A., Barnes, P., Tay, Y., Shazeer, N., Prabhakaran, V., Reif, E., Du, N., Hutchinson, B., Pope, R., Bradbury, J., Austin, J., Isard, M., Gur-Ari, G., Yin, P., Duke, T., Levskaya, A., Ghemawat, S., Dev, S., Michalewski, H., Garcia, X., Misra, V., Robinson, K., Fedus, L., Zhou, D., Ippolito, D., Luan, D., Lim, H., Zoph, B., Spiridonov, A., Sepassi, R., Dohan, D., Agrawal, S., Omernick, M., Dai, A.M., Pillai, T.S., Pellat, M., Lewkowycz, A., Moreira, E., Child, R., Polozov, O., Lee, K., Zhou, Z., Wang, X., Saeta, B., Diaz, M., Firat, O., Catasta, M., Wei, J., Meier-Hellstern, K., Eck, D., Dean, J., Petrov, S., Fiedel, N.: Palm: Scaling language modeling with pathways (2022) 4

22. Contributors, O.: Opencompass: A universal evaluation platform for foundation models. https://github.com/open-compass/opencompass (2023) 11
23. Contributors, X.: Xtuner: A toolkit for efficiently fine-tuning llm. https://github.com/InternLM/xtuner (2023) 11, 12, 23, 24
24. Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: Instructblip: Towards general-purpose vision-language models with instruction tuning (2023) 2, 4, 11, 12, 23, 24
25. DO CT OR OF, P.E.: MACHINE PERCEPTION OF THREE-DIMENSIONAL, SO LIDS. Ph.D. thesis, MASSACHUSETTS INSTITUTE OF TECHNOLOGY (1961) 2
26. Dong, X., Zhang, P., Zang, Y., Cao, Y., Wang, B., Ouyang, L., Wei, X., Zhang, S., Duan, H., Cao, M., Zhang, W., Li, Y., Yan, H., Gao, Y., Zhang, X., Li, W., Li, J., Chen, K., He, C., Zhang, X., Qiao, Y., Lin, D., Wang, J.: Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. arXiv preprint arXiv:2401.16420 (2024) 2, 11
27. Fang, Y., Wang, W., Xie, B., Sun, Q., Wu, L., Wang, X., Huang, T., Wang, X., Cao, Y.: Eva: Exploring the limits of masked visual representation learning at scale. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19358–19369 (2023) 4, 22
28. Fu, C., Chen, P., Shen, Y., Qin, Y., Zhang, M., Lin, X., Yang, J., Zheng, X., Li, K., Sun, X., et al.: Mme: A comprehensive evaluation benchmark for multimodal large language models. arXiv preprint arXiv:2306.13394 (2023) 4
29. Fu, S., Tamir, N., Sundaram, S., Chai, L., Zhang, R., Dekel, T., Isola, P.: Dreamsim: Learning new dimensions of human visual similarity using synthetic data (2023) 3, 10
30. Fu, X., Zhang, S., Kwon, G., Perera, P., Zhu, H., Zhang, Y., Li, A.H., Wang, W.Y., Wang, Z., Castelli, V., Ng, P., Roth, D., Xiang, B.: Generate then select: Open-ended visual question answering guided by world knowledge. In: Rogers, A., Boyd-Graber, J., Okazaki, N. (eds.) Findings of the Association for Computational Linguistics: ACL 2023. pp. 2333–2346. Association for Computational Linguistics, Toronto, Canada (Jul 2023). https://doi.org/10.18653/v1/2023.findings-acl.147, https://aclanthology.org/2023.findings-acl.147 4
31. Fu, X., Zhou, B., Chandratreya, I., Vondrick, C., Roth, D.: There's a time and place for reasoning beyond the image. In: Muresan, S., Nakov, P., Villavicencio, A. (eds.) Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 1138–1149. Association for Computational Linguistics, Dublin, Ireland (May 2022). https://doi.org/10.18653/v1/2022.acl-long.81, https://aclanthology.org/2022.acl-long.81 3, 4, 9
32. Fu, X., Zhou, B., Chen, S., Yatskar, M., Roth, D.: Interpretable by design visual question answering. arXiv preprint arXiv:2305.14882 (2023) 4
33. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2017) 4
34. Guan, T., Liu, F., Wu, X., Xian, R., Li, Z., Liu, X., Wang, X., Chen, L., Huang, F., Yacoob, Y., Manocha, D., Zhou, T.: Hallusionbench: An advanced diagnostic suite for entangled language hallucination & visual illusion in large vision-language models (2023) 5
35. Gupta, A., Dollar, P., Girshick, R.: Lvis: A dataset for large vocabulary instance segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 5356–5364 (2019) 3, 9

36. Harris, C., Stephens, M., et al.: A combined corner and edge detector. In: Alvey vision conference. vol. 15, pp. 10–5244. Citeseer (1988) 2
37. Hartley, R., Zisserman, A.: Multiple view geometry in computer vision. Cambridge university press (2003) 2
38. Hu, Y., Hua, H., Yang, Z., Shi, W., Smith, N.A., Luo, J.: Promptcap: Prompt-guided task-aware image captioning. arXiv preprint arXiv:2211.09699 (2022) 2, 4
39. Hu, Y., Liu, B., Kasai, J., Wang, Y., Ostendorf, M., Krishna, R., Smith, N.A.: Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. arXiv preprint arXiv:2303.11897 (2023) 4
40. Hu, Y., Stretcu, O., Lu, C.T., Viswanathan, K., Hata, K., Luo, E., Krishna, R., Fuxman, A.: Visual program distillation: Distilling tools and programmatic reasoning into vision-language models. arXiv preprint arXiv:2312.03052 (2023) 14
41. Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L.J., Shamma, D.A., et al.: Visual genome: Connecting language and vision using crowdsourced dense image annotations. International journal of computer vision **123**, 32–73 (2017) 3, 4
42. Lai, Z., Purushwalkam, S., Gupta, A.: The functional correspondence problem. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15772–15781 (2021) 3, 10
43. Li, B., Ge, Y., Ge, Y., Wang, G., Wang, R., Zhang, R., Shan, Y.: Seed-bench-2: Benchmarking multimodal large language models. arXiv preprint arXiv:2311.17092 (2023) 2, 5
44. Li, B., Wang, R., Wang, G., Ge, Y., Ge, Y., Shan, Y.: Seed-bench: Benchmarking multimodal llms with generative comprehension. arXiv preprint arXiv:2307.16125 (2023) 2, 3, 5
45. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023) 4, 23
46. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Computer Vision– ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13. pp. 740–755. Springer (2014) 3, 4, 5, 10
47. Liu, F., Emerson, G., Collier, N.: Visual spatial reasoning. Transactions of the Association for Computational Linguistics **11**, 635–651 (2023) 2, 9, 21
48. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning (2023) 4, 12, 23, 24
49. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning (2023) 11
50. Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., Lee, Y.J.: Llava-next: Improved reasoning, ocr, and world knowledge (January 2024), https://llava-vl.github.io/blog/2024-01-30-llava-next/ 2, 4, 8, 11, 12, 23, 24
51. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. Advances in neural information processing systems **36** (2024) 2, 11
52. Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., Chen, K., Lin, D.: Mmbench: Is your multi-modal model an all-around player? (2023) 2, 3, 5, 9, 13, 21
53. Liu, Y., Li, Z., Li, H., Yu, W., Huang, M., Peng, D., Liu, M., Chen, M., Li, C., Jin, L., et al.: On the hidden mystery of ocr in large multimodal models. arXiv preprint arXiv:2305.07895 (2023) 2
54. Liu, Y., Li, Z., Yang, B., Li, C., Yin, X., lin Liu, C., Jin, L., Bai, X.: On the hidden mystery of ocr in large multimodal models (2024) 4

55. Lowe, D.G.: Object recognition from local scale-invariant features. In: Proceedings of the seventh IEEE international conference on computer vision. vol. 2, pp. 1150–1157. Ieee (1999) 2

56. Lu, J., Clark, C., Lee, S., Zhang, Z., Khosla, S., Marten, R., Hoiem, D., Kembhavi, A.: Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action. arXiv preprint arXiv:2312.17172 (2023) 2

57. Lu, P., Bansal, H., Xia, T., Liu, J., Li, C., Hajishirzi, H., Cheng, H., Chang, K.W., Galley, M., Gao, J.: Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. arXiv preprint arXiv:2310.02255 (2023) 2, 5

58. Marr, D.: Vision: A computational investigation into the human representation and processing of visual information. MIT press (2010) 2

59. Marr, D., Poggio, T.: Cooperative computation of stereo disparity: A cooperative algorithm is derived for extracting disparity information from stereo image pairs. Science **194**(4262), 283–287 (1976) 2

60. Min, J., Lee, J., Ponce, J., Cho, M.: Spair-71k: A large-scale benchmark for semantic correspondence. arXiv preprint arXiv:1908.10543 (2019) 10

61. Minsky, M., Papert, S.: An introduction to computational geometry. Cambridge tiass., HIT **479**(480), 104 (1969) 2

62. OpenAI: Gpt-4 technical report (2023) 2, 4, 5, 8, 11, 12, 23, 24, 25

63. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952 (2023) 10

64. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) 4, 22, 23

65. Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., Komatsuzaki, A.: Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. arXiv preprint arXiv:2111.02114 (2021) 4

66. Schwenk, D., Khandelwal, A., Clark, C., Marino, K., Mottaghi, R.: A-okvqa: A benchmark for visual question answering using world knowledge. In: European Conference on Computer Vision. pp. 146–162. Springer (2022) 4

67. Shtedritski, A., Rupprecht, C., Vedaldi, A.: What does clip know about a red circle? visual prompt engineering for vlms. arXiv preprint arXiv:2304.06712 (2023) 14

68. Sun, J., Shen, Z., Wang, Y., Bao, H., Zhou, X.: LoFTR: Detector-free local feature matching with transformers. CVPR (2021) 14

69. Sun, Q., Fang, Y., Wu, L., Wang, X., Cao, Y.: Eva-clip: Improved training techniques for clip at scale. arXiv preprint arXiv:2303.15389 (2023) 4, 23

70. Tang, L., Jia, M., Wang, Q., Phoo, C.P., Hariharan, B.: Emergent correspondence from image diffusion. arXiv preprint arXiv:2306.03881 (2023) 14

71. Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A.M., Hauth, A., et al.: Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805 (2023) 2, 4, 8, 11, 12, 23, 24, 25

72. Team, I.: Internlm: A multilingual language model with progressively enhanced capabilities. https://github.com/InternLM/InternLM (2023) 12, 23, 24

73. Team, M.N.: Introducing mpt-7b: A new standard for open-source, commercially usable llms (2023), www.mosaicml.com/blog/mpt-7b, accessed: 2023-05-05 22

74. Torralba, A., Oliva, A.: Depth estimation from image structure. IEEE Transactions on pattern analysis and machine intelligence **24**(9), 1226–1238 (2002) 2

75. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al.: Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023) 4, 22
76. Wang, J.Y., Adelson, E.H.: Layered representation for motion analysis. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. pp. 361–366. IEEE (1993) 2
77. Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., Xu, J., Xu, B., Li, J., Dong, Y., Ding, M., Tang, J.: Cogvlm: Visual expert for pretrained language models (2023) 2, 11, 12, 23, 24
78. Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., Zhou, D.: Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171 (2022) 26
79. Wang, Z., Bao, J., Zhou, W., Wang, W., Hu, H., Chen, H., Li, H.: Dire for diffusion-generated image detection. arXiv preprint arXiv:2303.09295 (2023) 14
80. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., Zhou, D.: Chain-of-thought prompting elicits reasoning in large language models (2022) 4
81. Yan, A., Yang, Z., Wu, J., Zhu, W., Yang, J., Li, L., Lin, K., Wang, J., McAuley, J., Gao, J., et al.: List items one by one: A new data source and learning paradigm for multimodal llms. arXiv preprint arXiv:2404.16375 (2024) 4
82. Yang, J., Zhang, H., Li, F., Zou, X., Li, C., Gao, J.: Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. arXiv preprint arXiv:2310.11441 (2023) 14
83. Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., Zhao, H.: Depth anything: Unleashing the power of large-scale unlabeled data. In: CVPR (2024) 14
84. Yang, Z., Gan, Z., Wang, J., Hu, X., Lu, Y., Liu, Z., Wang, L.: An empirical study of gpt-3 for few-shot knowledge-based vqa. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 3081–3089 (2022) 2
85. Yang, Z., Li, L., Lin, K., Wang, J., Lin, C.C., Liu, Z., Wang, L.: The dawn of lmms: Preliminary explorations with gpt-4v (ision). arXiv preprint arXiv:2309.17421 **9** (2023) 9
86. Yu, W., Yang, Z., Li, L., Wang, J., Lin, K., Liu, Z., Wang, X., Wang, L.: Mm-vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint arXiv:2308.02490 (2023) 5
87. Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., Wei, C., Yu, B., Yuan, R., Sun, R., Yin, M., Zheng, B., Yang, Z., Liu, Y., Huang, W., Sun, H., Su, Y., Chen, W.: Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. arXiv preprint arXiv:2311.16502 (2023) 2, 3, 5, 9, 13
88. Ze, Y., Wang, X.: Category-level 6d object pose estimation in the wild: A semi-supervised learning approach and a new dataset. Advances in Neural Information Processing Systems **35**, 27469–27483 (2022) 3, 9
89. Zellers, R., Bisk, Y., Farhadi, A., Choi, Y.: From recognition to cognition: Visual commonsense reasoning. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019) 4
90. Zheng, L., Chiang, W.L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al.: Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in Neural Information Processing Systems **36** (2024) 4, 23

In the supplemental materials, Appendix A contains additional details on BLINK dataset collection and model inference, Appendix B provides more details of the baseline models, Appendix C includes experimental analyses on BLINK, and Appendix D discusses limitations.

## A    BLINK Details

### A.1    Visual Prompts Details

There are three types of visual prompts in BLINK: circles, boxes, and masks as shown in Figure 5. As for visual correspondence, functional correspondence, semantic correspondence, the red circles have radius 10px on images resized to 1024px height. For relative reflectance, we draw white circles to avoid color confusions. For object localization, the boxes are in red and green. For jigsaw, the masks are kept black. Since the examples in Figure 5 are different from the actual ones for illustrative purposes, we show some actual-sized example data as in Figures 9 to 19, with GPT-4V predictions attached.

```
gpt_query_template = (
"You are an AI assistant who will help me to match an answer with several options of a single-choice question. "
"You are provided with a question, several options, and an answer, and you need to find which option is most similar to the answer. "
"If the answer says things like refuse to answer, I'm sorry cannot help, etc., output (Z)"
"If the meaning of all options are significantly different from the answer, or the answer does not select any option, output (Z)"\
"Your should output one of the choices, (A),(B),(C),(D),(E) (if they are valid options), or (Z)\n"
"Example 1: \n"
"Question: Which point is closer to the camera?\nSelect from the following choices.\nOptions: (A) Point A\n(B) Point B\n(Z) Failed "
" Answer: Point B, where the child is sitting, is closer to the camera.\nYour output: (B)\n"
"Example 2: \n"
"Question: Which point is closer to the camera?\nSelect from the following choices.\nOptions: (A) Point A\n(B) Point B\n(Z) Failed "
" Answer: I'm sorry, but I can't assist with that request.\nYour output: (Z)\n"
"Example 3: \n"
"Question: Which point is corresponding to the reference point?\nSelect from the following choices.\nOptions: (A) Point A\n(B) Point B\n(Z) Failed "
" Answer:The reference point (REF) on the first image is at the tip of the pot, which is the part used to Poke if the pots were used for that action.
   Looking at the second image, we need to find the part of the object that would correspond to poking.\n(A) Point A is at the tip of the spoon's
handle, which is not used for poking.\n(B) Point B is at the bottom of the spoon, which is not used for poking.\n(C) Point C is on the side of the
pspoonot, which is not used for poking.\n(D) Point D is at the tip of the spoon, which is not used for poking.\n\nTherefore, there is no correct
answer in the choices\nYour output: (Z)\n"
"Example 4: \n"
"Question: {question}?\nOptions: {options}\n(Z) Failed\nAnswer: {prediction}\nYour output: ")
```

**Figure 8:** The evaluation prompts used for option label extraction.

### A.2    Spatial Relation Curation Process

We curate our samples from the Visual Spatial Reasoning [47] dataset. Each original sample contains an image and a claim, which is either true or false. One example being "`Caption: The cow is ahead of the person. Label: False.`" We reformat the claims into binary questions via GPT-3.5 [13], *e.g.* "`Question: Is the cow ahead of the person? Choices: (A) Yes (B) No Label: (B)`"

### A.3    Evaluation Prompts

Following MMBench [52], given model outputs, we first try to extract choices with exact matching (*e.g.*, for 'C', we try to match "C" and "(C)", etc). If failed,

we extract the choices using GPT-3.5 [13]. We provide GPT with the question, options, and model prediction, and then request GPT to align the prediction with one of the given options, or "Z", meaning that it fails to match to an option. Screenshot of the prompts we used are in Figure 8.

### A.4   Human Evaluation Protocol

We assign two humans (coauthors) for each task in BLINK and present their average scores as human performance. The human agreement scores range between 80-99%, with the lowest being on art style and functional correspondence, highest being relative depth, object localization, and forensics detection. Notice that the only exception is the IQ test score provided by two coauthors tested upon 100 sampled data, 50 for each, since it is hard to control or represent as average human performance.

### A.5   Dataset Statistics

Detailed statistics of BLINK are shown in Table 3.

| Statistics | Number |
|---|---|
| Total Questions | 3,807 |
| Total Images | 7,358 |
| Dev:Test | $1,901 : 1,906$ |
| Questions with Visual Prompts | 1,946 |
| Questions with Images (regions) as Choices | 2,747 |
| Questions with an Explanation | 300 |
| Questions with Multiple Images | 2,218 |
| * with 2 Images | 1,149 |
| * with 3 Images | 805 |
| * with 4 Images | 264 |

**Table 3:** Detailed statistics of the BLINK benchmark.

## B   Baseline Models

We evaluate BLINK on 16 various large multimodal LLMs. For most model families, we use the latest and best-performing available checkpoint to date. The list of baseline models are as follows: (i) MiniGPT-4-v2 [17] adapts EVA [27] as visual backbone, LLaMA2-chat (7B) [75] as language model backbone, and designs a linear projection layer for visual understanding abilities. (ii) OpenFlamingo [5] is an an open-source alternative to Flamingo [3] and we use the 9B checkpoint model, built upon CLIP [64] vision encoder and MPT-7B language model [73].

(iii - iv) InstructBLIP [24] uses CLIP [64] for vision encoder, and is fine-tuned based on BLIP-2 [45] with visual instruction data. We experiment with the 7B and 13B scales, both based on the Vicuna [90] language model for model scaling analysis. (v-x) We include various LLaVa [48] models from different sources for comparison: LLaVa-internLM2-7B which is fine-tuned upon InternLM2-Chat-7B [72] language model; LLaVa-v1.5-7B-xtuner and LLaVa-v1.5-13B-xtuner that are fine-tuned upon Vicuna [90] from xTuner [23]; and LLaVa-v1.5-7B, LLaVa-v1.5-13B, LLaVa-v1.6-34b from the original LLaVa papers [48,50]. Compared to the v1.5 checkpoints, v1.6 checkpoint uses more reasoning, OCR, and knowledge-enhanced training data. All of the LLaVa models build upon the CLIP [64] vision encoder. (vi-vii) Yi-VL-6B and Yi-VL-34B[4] are open-source models that have shown great performance on existing benchmarks. They use LLaVa structure with CLIP [64] encoder and connect with Yi-6B-Chat or Yi-34B-Chat language models[5]. (viii) CogVLM [77] adds a trainable visual expert module in the attention and FFN layers to bridge different modalities better. It uses EVA-CLIP [69] as vision encoder and Vicuna [90] as language backbone. (ix)Qwen-VL [7] includes several powerful models that show supreme performance on existing benchmarks. We use the best model checkpoint: Qwen-VL-MAX. (x) GeminiProVision [71] is one of the most powerful multimodal models, and we use the Gemini 1.0 Pro Vision version. (xi) Glaude 3 OPUS [1] is a recently released multimodal model that is tested to be state-of-the-art on various datasets. We use the most powerful version: OPUS, of the Claude 3 model family. (xii) GPT-4 [62] is known to be one of the most powerful multimodal models to date. We tested on three checkpoints: GPT-4V(ision), which is gpt-4-vision-preview; GPT-4 Turbo, which is gpt-4-turbo-2024-0409; and GPT-4o, which is gpt-4o-2024-05-13.

**GPT-4 Clarification**. Notice that the GPT-4 performances could change if the specific checkpoint gets updated. We tested GPT-4V(ision) in March 2024, and both of GPT-4 Turbo and GPT-4o in May 2024.

## C   Analysis

### C.1   Validation Set Results

We include detailed scores for each task on the validation set as in Table 4.

### C.2   How to deal with multiple-image inputs?

Among all the 16 baseline models, only 2 models: GPT-4V and Gemini Pro accept multi-image inputs. Other models, especially the open-source ones, only accept single-image inputs. Since 8 out of 14 of Blink tasks require multiple images input, a natural question is, how to deal with multiple-image inputs? To answer this question, we convert multiple images into concatenated single image, to analyze which format would achieve better performance on multi-image

---

[4] Model details can be found at https://huggingface.co/01-ai/Yi-VL-6B

[5] More details are at the official website at https://www.01.ai/

| | Validation (1,901) | Test (1,906) | Similarity (135) | Counting (120) | Depth (124) | Jigsaw (150) | Art (117) | Fun.Corr. (130) |
|---|---|---|---|---|---|---|---|---|
| Random Choice | 38.09 | 38.09 | 50 | 25 | 50 | 50 | 50 | 25 |
| Human | 95.67 | 95.70 | 96.70 | 93.75 | 99.19 | 99.00 | 95.30 | 80.77 |
| **Open-source multimodal LLMs** | | | | | | | | |
| MiniGPT-4-v2 [16] | 34.23 | 34.57 | 44.44 | 13.33 | 50.81 | 34.67 | 43.59 | 20.77 |
| OpenFlamingo-v2 [5] | 39.18 | 38.32 | 62.22 | 30.00 | 54.03 | 47.33 | 52.99 | 24.62 |
| InstructBLIP-7B [24] | 39.72 | 38.65 | 47.41 | 32.50 | 51.61 | 52.67 | 47.01 | 23.85 |
| InstructBLIP-13B [24] | 42.24 | 39.58 | 49.63 | 30.83 | 51.61 | 52.67 | 51.28 | 29.23 |
| LLaVA-internLM2-7B [72] | 37.71 | 36.06 | 48.89 | 55.00 | 57.26 | 28.67 | 29.06 | 23.85 |
| Yi-VL-6B [2] | 38.72 | 41.24 | 46.67 | 55.00 | 57.26 | 48.00 | 39.32 | 17.69 |
| Yi-VL-34B [2] | 41.68 | 42.78 | 51.11 | 52.50 | 50.00 | 52.67 | 45.30 | 31.54 |
| LLaVA-v1.5-7B-xtuner [23] | 39.36 | 40.81 | 47.41 | 45.83 | 51.61 | 52.67 | 47.01 | 20.00 |
| LLaVA-v1.5-13B-xtuner [23] | 42.00 | 41.31 | 47.41 | 48.33 | 54.03 | 52.00 | 47.01 | 30.00 |
| CogVLM [77] | 41.54 | 39.38 | 47.41 | 38.33 | 52.42 | 52.67 | 47.86 | 23.08 |
| LLaVA-v1.5-7B [48] | 37.13 | 38.01 | 47.41 | 40.00 | 52.42 | 11.33 | 47.01 | 20.00 |
| LLaVA-v1.5-13B [48] | 42.66 | 40.55 | 47.41 | 45.00 | 53.23 | 58.00 | 47.01 | 26.15 |
| LLaVA-v1.6-34B [50] | 46.80 | 45.05 | 48.89 | 66.67 | 67.74 | 54.67 | 43.59 | 20.77 |
| **API-based models** | | | | | | | | |
| Qwen-VL-Max [7] | 40.28 | 41.94 | 51.11 | 56.67 | 58.06 | 4.67 | 38.46 | 28.46 |
| Gemini Pro [71] | 45.16 | 45.72 | 52.59 | 52.50 | 40.32 | 57.33 | 50.43 | 24.62 |
| Claude 3 OPUS [1] | 44.05 | 44.11 | 72.59 | 47.58 | 32.67 | 65.81 | 21.54 |
| GPT-4V(ision) [62] | 51.14 | 51.26 | 78.52 | 60.83 | 59.68 | 70.00 | 79.49 | 26.15 |
| GPT-4 Turbo [62] | 54.61 | 53.89 | 80.74 | 57.50 | 66.13 | 69.33 | 79.49 | 24.62 |
| GPT-4o [62] | 60.04 | 59.03 | 72.59 | 49.17 | 74.19 | 55.33 | 82.91 | 40.77 |

| | Sem.Corr. (139) | Spatial (143) | Local. (122) | Vis.Corr. (172) | Multi-view (133) | Reflect. (134) | Forensic (132) | IQ (150) |
|---|---|---|---|---|---|---|---|---|
| Random Choice | 25 | 50 | 50 | 25 | 50 | 33.33 | 25 | 25 |
| Human | 96.07 | 98.25 | 98.00 | 99.42 | 92.48 | 95.14 | 100.00 | 80.00 |
| **Open-source multimodal LLMs** | | | | | | | | |
| MiniGPT-4-v2 [16] | 28.78 | 44.76 | 47.54 | 26.16 | 48.87 | 30.60 | 24.24 | 20.67 |
| OpenFlamingo-v2 [5] | 30.22 | 43.36 | 56.56 | 25.58 | 44.36 | 36.57 | 21.97 | 18.67 |
| InstructBLIP-7B [24] | 30.94 | 56.64 | 48.36 | 30.81 | 55.64 | 33.58 | 25.00 | 20.00 |
| InstructBLIP-13B [24] | 32.37 | 65.73 | 55.74 | 29.65 | 57.14 | 38.81 | 21.97 | 24.67 |
| LLaVA-internLM2-7B [72] | 27.34 | 76.22 | 50.00 | 27.91 | 44.36 | 32.09 | 5.30 | 22.00 |
| Yi-VL-6B [2] | 18.71 | 68.53 | 45.08 | 26.74 | 42.86 | 27.61 | 27.27 | 21.33 |
| Yi-VL-34B [2] | 19.42 | 71.33 | 51.64 | 26.74 | 44.36 | 38.81 | 23.48 | 24.67 |
| LLaVA-v1.5-7B-xtuner [23] | 28.78 | 68.53 | 36.89 | 29.07 | 38.35 | 29.85 | 36.36 | 18.67 |
| LLaVA-v1.5-13B-xtuner [23] | 30.94 | 69.93 | 45.08 | 29.65 | 44.36 | 38.81 | 25.76 | 24.67 |
| CogVLM [77] | 33.09 | 63.64 | 52.46 | 29.65 | 54.14 | 29.85 | 30.30 | 26.67 |
| LLaVA-v1.5-7B [48] | 23.02 | 61.54 | 56.56 | 25.58 | 51.88 | 39.55 | 23.48 | 20.00 |
| LLaVA-v1.5-13B [48] | 32.37 | 67.83 | 52.46 | 29.07 | 44.36 | 36.57 | 31.82 | 26.00 |
| LLaVA-v1.6-34B [50] | 23.74 | 74.83 | 59.02 | 30.81 | 62.41 | 31.34 | 44.70 | 26.00 |
| **API-based models** | | | | | | | | |
| Qwen-VL-Max [7] | 23.02 | 69.93 | 48.36 | 31.40 | 51.88 | 36.57 | 43.94 | 21.33 |
| Gemini Pro [71] | 26.62 | 74.83 | 53.28 | 42.44 | 44.36 | 38.81 | 50.76 | 23.33 |
| Claude 3 OPUS [1] | 25.18 | 58.04 | 51.64 | 36.63 | 56.39 | 26.87 | 46.21 | 24.67 |
| GPT-4V(ision) [62] | 28.78 | 72.73 | 54.92 | 33.72 | 55.64 | 38.81 | 34.09 | 22.67 |
| GPT-4 Turbo [62] | 30.94 | 69.23 | 52.46 | 52.33 | 52.63 | 32.84 | 63.64 | 32.67 |
| GPT-4o [62] | 53.96 | 69.23 | 59.84 | 75.00 | 59.40 | 37.31 | 79.55 | 31.33 |

**Table 4: Results of different models on the BLINK validation set**. The first row shows task names and number of instances.

understanding. Specifically, we place the images horizontally, with a black margin in between. We evaluate GPT-4V and Gemini Pro with concatenated images and show results in Table 5.

From the experiment results, GPT-4V has shown a consistent decline in performance across all tasks when taking concatenated images as input, with the

|  | Similarity | Jigsaw | Art | Fun.Corr. | Sem.Corr. | Vis.Corr. | Multi-view | Forensic |
|---|---|---|---|---|---|---|---|---|
| Random Choice | 50 | 50 | 50 | 25 | 25 | 25 | 50 | 25 |
| Human | 96.70 | 99.00 | 95.30 | 80.77 | 96.07 | 99.42 | 92.48 | 100.00 |
| Gemini Pro [71] | 55.88 | 54.00 | 49.57 | **32.31** | 22.14 | **37.21** | 41.35 | **45.45** |
| * concatenate images | 42.65 | 45.33 | 48.72 | 30.77 | 27.86 | 23.84 | 41.35 | 36.36 |
| GPT-4V(ision) [62] | **83.09** | **62.67** | **78.63** | 31.54 | **30.00** | **37.21** | **58.65** | 30.30 |
| * concatenate images | 71.32 | 57.33 | 67.52 | 22.31 | 22.86 | 25.00 | 57.89 | 25.00 |

**Table 5: Effect of concatenating multiple images on the BLINK val set.**

biggest decrease in jigsaw and least decrease in multi-view reasoning. However, the impact of concatenating images to Gemini Pro is task-dependent, with the performance decreasing in most tasks while increasing in semantic correspondence and remaining the same in multi-view reasoning.

### C.3    Error analysis

**Open-source multimodal LLMs make similar errors.** Our comparative analysis of diverse multimodal LLMs reveals striking similarities in the cases where they fail at, highlighting that these shared mistakes are largely influenced by their vision encoder, rather than differences in model size or language model components. This is particularly apparent in the comparison between LLaVA-v1.5-7B (1187 mistakes in total) and LLaVA-v1.5-13B (1147 mistakes in total), two models of different sizes that nonetheless demonstrated 899 common mistakes. In a similar vein, when we compared LLaVA-v1.5-7B with other equal-sized models using different language model components, like LLaVA-internLM2-7B, the number of common errors remained high (959 mistakes). Whereas LLaVA-v1.5-7B only shares 782 and 655 common mistakes with QwenVLMax and GPT4V, respectively.

**GPT-4V Errors:** For each task, 10 error instances were randomly selected, and we manually analyze the total of 140 error instances sampled randomly across all tasks as follows: Recognition failure on detailed small regions or edges (28.5%) : the model fails to tell nuanced details, especially circles in visual correspondence, semantic correspondence, functional correspondence, relative depth, relative reflectance and boxes in object localization; Failure to detect the location of the circled point(20%): the model fails to locate the circled point labeled in the images; Failure to recognize spatial relations (14.3%): the model fails to identify the spatial relations between left and right, or up and down; Reasoning errors (12.9%): while the model correctly interprets the images and the question, it fails to derive accurate reasoning for inference; Failure to convey the overall scene impression (8.6%): the model fails to adequately capture the general atmosphere or setting of a scene; Rejection to answer (6.4%): the model refuses to generate an answer; Failure to ground or infer items mentioned in the question (5.7%): The model is unable to locate the specific item referenced in the question within the image.

### C.4    Does self-consistency help?

To verify whether self-consistency [78] will improve the performance on BLINK, we conduct five runs of GPT-4V with temperature set to 1.0. The self-consistency score is 48.38%, averaged across all task, and the average single run score is 38.15%. Note that the performance for temperature=1 decreases much compared to our default temperature=0 setting, where greedy decoding is used. From our observation, this decrease is because GPT-4V tends to face firewalls and reply "Sorry I cannot help/assist with the question" when temperature is bigger than 0.

## D    Limitations

BLINK makes use of data from existing image datasets, and does not cover all the visual perception abilities in the wild. For the forensics detection task, we manually collected images that are publicly available from online search. We have made every effort to ensure that the images included in this paper are used in accordance with applicable copyright laws and are properly credited. However, if you are the copyright owner of any image included in our work and believe that its use conflicts with your licensing agreements, please contact us directly. We are committed to addressing any legitimate concerns promptly.

**Visual correspondence**

**Prompt**:
A point is circled on the first image, labeled with REF. We change the camera position or lighting and shoot the second image. You are given multiple red-circled points on the second image, choices of "A, B, C, D" are drawn beside each circle.
Which point on the second image corresponds to the point in the first image? Select from the following options.
(A) Point A
(B) Point B
(C) Point C
(D) Point D

**Gold Answer**: C

**GPT-4V**:
To determine which point on the second image corresponds to the point labeled "REF" on the first image, we need to look for distinctive patterns or colors around the reference point that can be matched in both images.

In the first image, the "REF" point is located on a dark blue square. To the left of this square is a light blue square, and above it is a green square. To the right of the green square is a yellow square, and to the right of the yellow square is a red square.

Now, looking at the second image, we need to find a similar arrangement of colored squares. Point A is on a red square, which does not match the dark blue square of the reference point. Point B is on a blue square, but the surrounding colors do not match the pattern we're looking for. Point C is on a green square, which also does not match the dark blue reference square. Point D, however, is on a dark blue square, and to the left of it is a light blue square, above it is a green square, to the right of the green square is a yellow square, and to the right of the yellow square is a red square.

This pattern matches the arrangement of colors around the "REF" point in the first image. Therefore, the point on the second image that corresponds to the "REF" point in the first image is:

(D) Point D

**Figure 9:** Examples of actual-sized data in BLINK with GPT-4V predictions.(1/11)

Relative reflectance



**Prompt**:
Two points are annotated on the image, labeled by A and B. Consider the surface color of the points (the albedo of the surface, without the effect of shading). Which point has darker surface color, or the colors is about the same? Select from the following choices.
(A) A is darker
(B) B is darker
(C) About the same

**Gold Answer**: B

**GPT-4V**:
In this image:
(A) A is darker
Point A is on a pillow and appears slightly shadowed, while point B is on the headboard and is well-lit. Even considering the albedo and ignoring the effects of shading, the pillow's fabric at point A is likely to be a bit darker than the wooden headboard at point B.

Relative depth



**Prompt**:
Two points are circled on the image, labeled by A and B beside each circle. Which point is closer to the camera? Select from the following choices.
(A) A is closer
(B) B is closer

**Gold Answer**: B

**GPT-4V**:
Point A is closer to the camera. It is located on the edge of the roof closer to the front of the building and thus closer to the perspective from which the photo was taken. Point B is further back on the roofline.

**Figure 10:** Examples of actual-sized data in BLINK with GPT-4V predictions. (2/11)

**Spatial reasoning**

**Prompt**:
Is the laptop touching the teddy bear? Select from the following choices.
(A) Yes
(B) No

**Gold Answer**: A

**GPT-4V**:
The laptop does not appear to be touching the teddy bear. There seems to be a small gap between the bear and the laptop. So the answer would be:

(B) no



**Multi-view reasoning**

**Prompt**:
The images are frames from a video. The video is shooting a static scene. The camera is either moving clockwise (left) or counter-clockwise (right) around the object. The first image is from the beginning of the video and the second image is from the end. Is the camera moving left or right when shooting the video? Select from the following options.
(A) left
(B) right

**Gold Answer**: A

**GPT-4V**:
The camera is moving to the right (clockwise) around the object. This is evident because the second image shows more of the left side of the product packaging compared to the first image, indicating that the camera has moved from the left to the right around the object.

**Figure 11:** Examples of actual-sized data in BLINK with GPT-4V predictions. (3/11)

**Prompt:**
Given the first image with the lower right corner missing, can you tell which one of the second image or the third image is the missing part? Imagine which image would be more appropriate to place in the missing spot. You can also carefully observe and compare the edges of the images. Select from the following choices.
(A) the second image
(B) the third image

**Gold Answer**: B

**GPT-4V**:
Based on the visual information provided, the correct choice to fill in the missing lower right corner of the first image would be:
(A) second image
The second image shows a continuation of the wall with similar patterns of damage and discoloration that match the wall in the first image. The edges of the second image also seem to align with the exposed edges of the first image, suggesting that it is indeed the missing part of the first image. The third image, on the other hand, shows people and does not appear to match the wall segment in the first image.

**Figure 12:** Examples of actual-sized data in BLINK with GPT-4V predictions. (4/11)

Art style

(A)          (B)

**Prompt**:
Some most common art painting styles include Realism, Impressionism, Expressionism, Pop Art, and Cubism. Given the following images of art paintings, use the first image as the reference image, and determine which one of the second or the third image shares the same style as the reference image? Select from the following choices.
(A) the second image
(B) the third image

**Gold Answer**: B

**GPT-4V**:
The first image is a sketch that appears to be in a style that is more illustrative or expressionistic, given the exaggerated features and the somewhat abstract representation of the figure. It does not neatly fit into any of the common art styles you've mentioned like Realism, Impressionism, Pop Art, or Cubism.

The second image is a painting that depicts a scene with two figures and is rendered in a style that aims to represent the scene realistically, with attention to detail, light, and shadow. This style is closest to Realism.

The third image is an abstract composition that simplifies and fragments the forms into geometric shapes, which is characteristic of Cubism.

Given these observations, none of the images share the exact style of the reference image. However, if we must choose the one that is stylistically closest, we would select:
(A) the second image

This is because the second image, while realistic, has a more narrative and possibly expressive quality.

**Figure 13:** Examples of actual-sized data in BLINK with GPT-4V predictions. (5/11)

**Object localization**

Target: Bus

**Prompt:**
A bounding box is an annotated rectangle surrounding an object. The edges of bounding boxes should touch the outermost pixels of the object that is being labeled.
Given the two bounding boxes on the image, labeled by A and B, which bounding box more accurately localizes and encloses the bus? Select from the following options.
(A) Box A
(B) Box B

**Gold Answer:** A

**GPT-4V**: Box A more accurately localizes and encloses the bus (vehicle). It fits closely to the edges of the bus, touching the outermost pixels of the object as per the definition of a bounding box provided. Box B, on the other hand, is around a person on a bicycle and does not enclose the bus at all.



**Counting**

**Prompt**:
How many tall buildings have a flat top? Select from the following choices.
(A) 2
(B) 1
(C) 0
(D) 3

**Gold Answer**: C

**GPT-4V**: The image shows a clock tower and some structures behind it, but it's not clear if the structures are tall buildings or part of the same building complex. Additionally, the resolution and angle of the photo do not provide enough detail to definitively determine the shape of the tops of the buildings in the background. Therefore, I cannot accurately answer the question about the number of tall buildings with a flat top based on this image.

**Figure 14:** Examples of actual-sized data in Blink with GPT-4V predictions. (6/11)

Forensic detection

**Prompt:** You are a judge in a photography competition, and now you are given the four images. Please examine the details and tell which one of them is most likely to be a real photograph. Select from the following choices.
(A) the first image
(B) the second image
(C) the third image
(D) the fourth image

**Gold Answer:** A

**GPT-4V**:
Based on the examination of the four images provided, I would say that:
(B) the second image is most likely to be a real photograph.
The reasons for this assessment are as follows:
- The lighting in the second image appears natural and soft, with subtle shadows that suggest a natural light source.
- The apples have a more detailed texture, with slight imperfections and variations in color that are typical of real fruit.
- The composition seems less uniform and more random, which is more likely in a natural setting as opposed to a digitally created one.

The other images, while they may also be photographs, have qualities that could suggest they have been digitally manipulated or are less natural in their presentation. For example, the uniformity and saturation of colors in the first, third, and fourth images might be less typical of a straightforward photograph without any digital enhancements.

**Figure 15:** Examples of actual-sized data in BLINK with GPT-4V predictions. (7/11)

**Figure 16:** Examples of actual-sized data in BLINK with GPT-4V predictions. (8/11)

**Visual similarity**

(A)  (B)

**Prompt:**
Given three similar but different images, take the first image as reference. Can you tell which one of the latter two images is most similar to the first one? Select from the following choices.
(A) the second image
(B) the third image

**Gold Answer:** A

**GPT-4V**:
To determine which image is most similar to the first one, let's compare the key features of the waterfall in the first image with the other two:

1. The first image shows a waterfall with a wide water flow, a significant amount of mist at the bottom, and a lush green environment surrounding it. The sunlight is visible at the top, suggesting the time of day might be either morning or late afternoon.

2. The second image shows a waterfall with a narrower flow and less mist at the bottom. The surrounding rock formations are more prominently visible, and the environment appears to be less lush than in the first image.

3. The third image shows a waterfall with a wider flow, similar to the first image. There is a visible mist at the bottom, and the surrounding environment is lush with greenery. The color of the water and the presence of foliage are more similar to the first image.

Based on these observations, the third image (B) is most similar to the first

**Figure 17:** Examples of actual-sized data in BLINK with GPT-4V predictions. (9/11)

**Prompt**:
Humans can find corresponding points for different objects in the same category. For instance, if there are images of two different cats, then the left ear tip of one cat corresponds to the left ear tip of the other cat, and the right front paw of one cat corresponds to the right front paw of the other cat. Given the following two images, a reference point is annotated on the first image, labeled with REF. You are given multiple red-circled points on the second image, choices of "A, B, C, D" are drawn beside each circle. Select between the choices on the second image and find the corresponding point for the reference point. Which point is corresponding to the reference point?
Select from the following choices.
 (A) Point A
(B) Point B
(C) Point C
(D) Point D

**Gold Answer**: D

**GPT-4V**:
The reference point labeled "REF" on the first image is located on the cat's nose. To find the corresponding point on the second image, we need to locate the nose of the cat in that image. The nose of the cat in the second image is circled and labeled with "C." Therefore, the correct corresponding point for the reference point is:

(C) Point C

**Figure 18:** Examples of actual-sized data in BLINK with GPT-4V predictions. (10/11)

Functional correspondence

**Prompt**: Humans can find corresponding points for the same action between different objects. For instance, if a person uses a pot versus a hammer to "Mash Pound", then the handle of the pot will be the corresponding point to the handle of the hammer because they serve the same function for the action -- to hold; and the bottom of the pot will be the corresponding point to the face of the hammer because they both mash the other object. Given the following two images, a reference point is annotated on the first image, labeled with REF. You are given multiple red-circled points on the right image, choices of "A, B, C, D" are drawn beside each circle. Select from the choices on the second image and find the corresponding point for the reference point, if we use both items for the action: "Mash/Pound". Which point is corresponding to the reference point? Select from the following choices.
(A) Point A
(B) Point B
(C) Point C
(D) Point D

**Gold Answer**: D

**GPT-4V**:
The reference point (REF) on the first image is located at the bottom of the shoe, which is the part that would be used to "Mash/Pound" something by stepping or stomping on it.

(A) Point A is at the bottom of the bottle, which would be the part that could be used to mash or pound something if you were to use the bottle in an upside-down position.
(B) Point B is on the body of the bottle, which would not be effective for mashing or pounding.
(C) Point C is on the neck of the bottle, which is also not suitable for mashing or pounding.
(D) Point D is on the cap of the bottle, which is not the part that would correspond to the bottom of the shoe for the action of mashing or pounding.

Therefore, the correct choice is: (A) Point A

**Figure 19:** Examples of actual-sized data in BLINK with GPT-4V predictions. (11/11)