# VisionArena: 230K Real World User-VLM Conversations with Preference Labels

Christopher Chou*
Stanford

Lisa Dunlap*
UC Berkeley

Koki Mashita
UC Berkeley

Krishna Mandal
UC Berkeley

Trevor Darrell
UC Berkeley

Ion Stoica
UC Berkeley

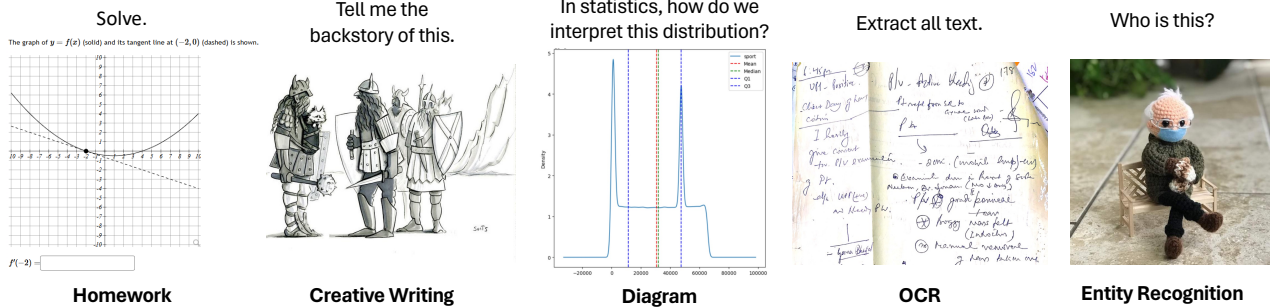Joseph E. Gonzalez
UC Berkeley

Wei-Lin Chiang
UC Berkeley

Figure 1. **Samples from VisionArena Conversations.** VisionArena contains conversations from real users covering a variety of domains.

## Abstract

*With the growing adoption and capabilities of vision-language models (VLMs) comes the need for benchmarks that capture authentic user-VLM interactions. In response, we create VisionArena, a dataset of 230K real-world conversations between users and VLMs. Collected from Chatbot Arena — an open-source platform where users interact with VLMs and submit preference votes — VisionArena spans 73K unique users, 45 VLMs, and 138 languages. Our dataset contains three subsets:* **VisionArena-Chat**, *200k single and multi-turn conversations between a user and a VLM;* **VisionArena-Battle**, *30K conversations comparing two anonymous VLMs with user preference votes; and* **VisionArena-Bench**, *an automatic benchmark of 500 diverse user prompts that efficiently approximate the live Chatbot Arena model rankings. Additionally, we highlight the types of question asked by users, the influence of response style on preference, and areas where models often fail. We find open-ended tasks like captioning and humor are highly style-dependent, and current VLMs struggle with spatial reasoning and planning tasks. Lastly, we show fine-tuning the same base model on VisionArena-Chat outperforms Llava-Instruct-158K, with a 17-point gain on MMMU and a 46-point gain on the WildVision benchmark. Dataset at* https://huggingface.co/lmarena-ai.

---
*Equal contribution.

## 1. Introduction

Visual language models (VLMs) [2, 3, 33, 37] are being increasingly used in a wide range of real-world applications including image captioning and story telling, document understanding, web development, and embodied systems. While these models have made remarkable progress on a wide range of benchmarks [4, 14, 18, 19, 44, 47], existing VLM benchmarks focus largely on static, single-turn tasks with predetermined correct answers, overlooking the open-ended, evolving nature of real-world user interactions. They also rarely capture multi-turn dialogue, incorporate diverse context, or reflect the fluidity of user intent. As such, they provide a simplified snapshot of VLM capabilities.

Understanding these real-world interactions across a variety of tasks is essential for developing models that align with human expectations and perform effectively. To address this, previous works such as Chatbot Arena [9] and WildVision [25, 26] crowdsource evaluation by hosting platforms where users can freely interact with pairs of VLMs and provide preference votes. Building off of these works, we introduce **VisionArena**, a dataset of 230K real-world conversations between users and 38 VLMs in 135 languages, collected through the Chatbot Arena platform. VisionArena consists of:

- **VisionArena-Chat**: 200,000 single and multi-turn chat logs between users and VLMs, spanning 138 languages, 73k users, and 45 open source and proprietary VLMs.

| Dataset | # Convs | # Models | # Users | # Langs | % Unique Images per Sample | Avg. # Turns | Avg. # Tokens per Prompt | Avg. # Tokens per Response | Human Preference |
|---|---|---|---|---|---|---|---|---|---|
| LMSYS-Chat-1M | 1,000,000 | 25 | 13,500 | 35 | - | 2.0 | 36.9 | 214.2 | No |
| WildVision-Battle | 10,383 | 19 | - | 28 | 56.2 | 1.2 | 57.8 | 131.7 | Yes |
| WildVision-Chat | 45,170 | 9 | - | 26 | 33.4 | 1.4 | 81.3 | 171.6 | No |
| VisionArena-Battle | 30,000 | 17 | 14,031 | 90 | 76.4 | 1.3 | 90.2 | 393.6 | Yes |
| VisionArena-Chat | 200,000 | 45 | 72.933 | 138 | 62.1 | 1.5 | 184.1 | 634.3 | No |

Table 1. **Dataset Comparison.** Compared to previous VLM preference benchmarks, VisionArena contains 3x the amount of data, with more users, language, models, unique images, and conversation turns.

- **VisionArena-Battle**: 30,000 conversations where users interact with two anonymized VLMs, along with preference votes indicating which response they prefer.
- **VisionArena-Bench**: An automatic benchmark consisting of 500 diverse user prompts that can be used to cheaply approximate model rankings via automatic benchmarking with VLM as a judge.

We conduct analysis of these datasets and construct a set of popular question categories including captioning, OCR, humor, creative writing, entity recognition, and diagram understanding. We also explore the influence of stylistic properties of responses such as response length, markdown, and specificity on human preference. We find that more open-ended questions like captioning and humor are heavily influenced by style, which causes certain models like InternVL to have a disproportionately higher ranking in these categories. We provide this metadata with VisionArena to enable further analysis. We also highlight common failure modes of VLMs and provide a small curated set of user prompts where top proprietary models fail, including complex spatial reasoning and planning tasks.

Next, we demonstrate how VisionArena can be used to improve VLMs through instruction finetuning. Compared to LLaVA-Instruct-158K [24], by finetuning on data from VisionArena-Chat, models show a 17 point improvement in MMMU [44] and a 46 point improvement on the human preference benchmark WV-Bench [26]. In addition to VisionArena, we also release this finetuned model. Lastly, we build on existing work in automatic benchmarking of VLMs to show that evaluating on the 500 prompts in VisionArena-Bench results in a model lineup that is consistent with the much larger online preference leaderboard Chatbot Arena. When compared with other automatic preference benchmarks like WildVision-Bench, VisionArena-Bench is far more predictive of the online Chatbot Arena VLM leaderboard performance, which contains over 100,000 user votes as of October 23rd, 2024. We believe that VisionArena is a valuable resource to better understand how people are currently using VLMs and will be the foundation for research in VLM development and evaluation. In the future we plan to continue regular data releases including a large variety of models and multi-image conversations.

## 2. Related Works

**Crowdsourced Evaluations.** In the past few years, several platforms have emerged that aim to crowdsource evaluation for LLMs and VLMs by allowing users to provide preference votes. These platforms, such as Chatbot Arena [9], allow anyone to freely engage in open-ended conversations with state-of-the-art commercial and open-source models. Users are able to directly chat with specific models or chat with pairs of anonymous models side-by-side. In the anonymous side-by-side mode, users can provide direct feedback on which responses they preferred, which is used to build a leaderboard. WildVision [26] adopts a similar style to Chatbot Arena except that users interact with VLMs instead of LLMs. Our platform builds upon these works by creating a unified interface that allows users to chat with *either* LLMs or VLMs.

**Public Chat Datasets:** LMSYS-Chat-1M [48], OpenAssistant [21], and WildChat-1M [10, 21] are all public datasets constructed by capturing users conversations with state-of-the-art LLMs. These datasets have been highly influential because they represent more natural human conversations and often contain reward feedback signals. However, because these datasets capture text-only conversations, they do not provide insight into how users incorporate images into conversations and the behavior of VLMs. Building on the success of public chat datasets, there is a recent effort to extend the public chat datasets to the visual domain with WildVision, a dataset of 45k chat logs with 9 visual question answering models and 10.4K battle logs across 19 models. In contrast, our VisionArena data set contains 200K chat logs and 30K battle logs across 40+ models, including all of the strong proprietary models and many open-source models, making it the largest and most complete VLM conversation dataset to date. See Table 1 for dataset comparison.

**VLM Benchmarks:** Currently, VLM benchmarks are typically static datasets that have close-ended ground truth answers (either multiple-choice or predefined-string). Some popular examples of these benchmarks include MMMU [44], DocVQA [29], MME [12], and VQA 2.0 [15]. To combat against static nature of datasets and minimize test-set contamination, live benchmarks are also available. For

example, LiveXiv incorporates updated ArXiv manuscripts for VQA [38]. To incorporate benchmarking on open-ended responses, there has been a trend towards using strong models (e.g. GPT-4o) for VLM-as-a-judge to approximate human preference. Some notable examples include WildVision-Bench [26] and Prometheus-Vision [22]. We similarly adopt the VLM-as-a-judge framework to create VisionArena-Bench, curated from questions from Vision-Arena, allowing it to be crowdsourced, open, and live. We show that using VisionArena-Bench, we achieve better correlation and agreement with the VLM leaderboard on Chatbot Arena, which itself has 100x more votes.

## 3. Dataset and Platform

### 3.1. Interface

VisionArena-Chat and VisionArena-Battle were collected from Chatbot Arena [9], an open-source platform for evaluating large language models by human preference. On our platform, users are able to directly chat with specific models (direct chat) or chat with pairs of anonymous models side-by-side (battle mode). In battle mode, users can provide direct feedback on which responses they preferred, which is used to build a leaderboard. We refer to these anonymous side-by-side chats as 'battles'. An example of our interface can be found in the supplemental. Unlike previous VLM crowdsourcing platforms, we integrate LLMs and VLMs into one unified chat interface with a simple routing mechanism. In side-by-side chat, if a user uploads an image in the first turn of their conversation, we automatically select from two available VLMs; otherwise we sample from the available LLMs. We then collect the votes from the image conversations to compute the VLM leaderboard. We believe this encourages users who may have initially been interested in interacting with LLMs to interact with VLMs as well. In Sec. 4, we show that our conversations do have important distributional differences with WildVision.

Before using our service, users must accept terms of use, giving us their consent to store and release the conversation data. The platform is free to use, and there is no registration process. We are supported by sponsorships with inference providers. VisionArena is a subset of conversations collected from February 2024 to September 2024. Given the language and question distribution of our collected conversations (Fig. 15, Fig. 5), the majority of our users are likely located in North America, Europe, and East Asia and work in STEM related fields.

To encourage user interaction, we provide a 'random image' button, which samples from a preset bank of images from 5 datasets: NewYorker [20], ChartVQA [28], DocVQA [29], TextVQA [39], and WikiArt [41]. We exclude these in VisionArena-Battle as we aim to capture the natural distribution of user inputs when computing leaderboard rankings, but we do include conversations with preset images in VisionArena-Chat, which make up around 15% of conversations.

**Moderation.** We apply several moderation steps before sending the prompt to the model provider and perform data cleaning procedures before releasing this data to the public. Before the user receives the response from the model, we perform (not safe for work) NSFW and (child sexual abuse material) CSAM [31] image detection and then tag and terminate conversations that contain sexual, hateful, or violent content. For battles, we also perform OpenAI text moderation [34] on user text prompts and discard any responses which contain a violation. For direct chats, we only perform OpenAI text moderation on proprietary models to follow their usage policies. We do not perform text moderation on prompts for open-source models, which opens this data for future analysis.

Finally, as part of our data release process we use Google's Vision API [13] to remove personally identifiable information (PII) from both images and text, removing any content containing human faces or identifiable details. However, these automated detectors are not infallible, so our dataset may still contain NSFW content or PII. We encourage users who find such instances to notify the authors so the material can be removed.

### 3.2. From Preference to Leaderboard Ranking

Using preference votes from pairwise battles in anonymous side-by-side chat, we apply a Bradley-Terry (BT) model [6] to estimate the relative strengths of models through logistic regression. The model's coefficients serve as *arena scores*, which determine the leaderboard rankings.

Let $n$ denote the number of pairwise comparisons (battles) and $M$ the number of models. For each battle $i \in [n]$, we define:

- $X_i \in \mathbb{R}^M$: $X_{i,m} = 1$ if model $m$ is presented first to the judge, $X_{i,m} = -1$ if presented last, and 0 otherwise.
- $Y_i \in 0, 1$: The outcome, where 1 indicates the first model won.

The BT model estimates model strengths $\beta \in \mathbb{R}^M$ through logistic regression:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^M} \frac{1}{n} \sum_{i=1}^{n} \text{CE}(\sigma(X_i^\top \beta), Y_i) \qquad (1)$$

where CE represents the cross-entropy loss and $\sigma$ is the sigmoid function. The BT coefficients $\hat{\beta}$ are the ratings associated with each of the VLMs in the arena. These BT ratings are used to create the ordered ranking of models on the leaderboard. We bootstrap the BT rating estimate 100 times to construct a confidence interval for each rating (Fig. 2).
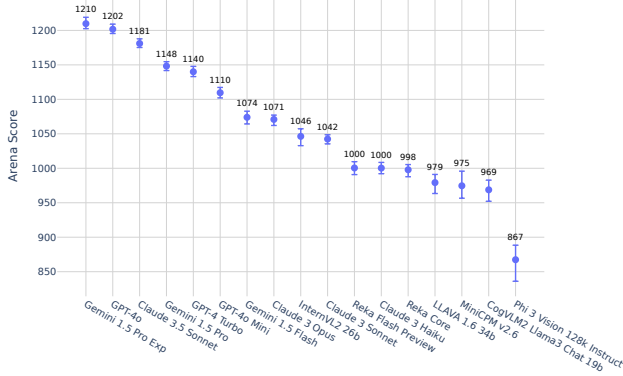
Figure 2. **Bootstrap B.T. model scores for VisionArena-Battle.** Proprietary models like Gemini 1.5 Pro and GPT-4o are at the top of the leaderboard, with open models like Llava 1.6, MiniCPM, CogVLMv2, and Phi3 obtaining the lowest ratings. InternVL2 is the highest rated open model, although as shown in Section 4.4, this is largely due to response style rather than model capability.
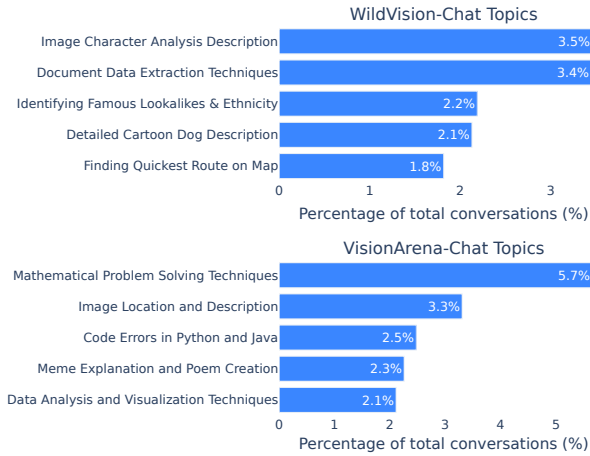


Figure 3. **Comparison of top 5 topic clusters between WildVision-Chat and VisionArena-Chat.** Compared to WildVision, the most popular topics clusters in VisionArena capture more real world tasks, specifically in STEM fields.

## 4. Data Analysis

In the following section, we (1) compare user votes to expert annotators (2), analyze and categorize the distribution of user conversations (3), compute per category leaderboard (4), measure the impact of response style on human preference, and (5) provide examples of difficult questions.

### 4.1. Comparing with Experts and the VisionArena

To inspect the quality of the human preference votes on our platform, we check their alignment with experts opinions. We sample 5 battles in English for each model pair and have 4 experts (PhD Students) label the battles based on their preference, creating an expert-labeled dataset of 516 responses.

We then compute BT scores using the expert labels and compare those scores with the BT scores computed directly on the VisionArena. We then compute the Pearson correlation coefficient [35] and the Spearman rank correlation coefficient [40]. We obtain a Pearson correlation of 0.88 indicating a strong linear predictive relationship between the BT scores computed by experts and those obtained from the live VisionArena. We obtain a Spearman rank correlation of 0.87 indicating high agreement in the ordering between the leaderboard rankings and the ranking obtained by our expert labelers on the small subset of data.

Additionally, for a subset of 100 battles labeled by 3 experts, we observe an agreement of 0.72 (excluding ties) and 0.56 (including ties) between users and expert annotators, compared to 0.77 (excluding ties) and 0.59 (including ties) among expert annotators themselves, further demonstrating the reliability of user votes.

### 4.2. What types of questions do people ask?

We perform topic modeling analysis on the VisionArena-Chat prompts. Following the BERTopic framework, We randomly sample 50K English conversations and embed the documents using CLIP-ViT-B-32, perform DBSCAN clustering, and use GPT-4o to summarize each cluster [16]. We plot our top 5 clusters in Fig. 3 (the top 20 clusters can be found in the supplement). We find that many people use VLMs to solve math and code problems, identify paintings and geographical locations, perform data analysis on tables and diagrams, explain humorous images, and create stories based on images. Notably, VisionArena-Chat contains important use cases not seen in WildVision-Chat including coding and web UI design problems, handwritten text extraction, and diagram analysis. Manually inspecting the clusters, we also see that WildVision-Chat's clusters are often very specific to a certain task (e.g. "Detailed Cartoon Dog Description", "Rice Leaf Disease Identification"), while VisionArena-Chat's cluster descriptions are broader which indicates the diversity of our data. Surprisingly, the majority of our questions require OCR, and we receive a large number of homework problems and diagram understanding questions. In the following section, we construct categories for each of these major use cases.

### 4.3. Prompt Categories

Based on the clustering analysis and manual inspection, we manually define 8 non-disjoint categories that reflect the vast majority of prompts and test different capabilities of the VLM. We use Gemini 1.5 Flash to classify each prompt, using both the image and text, into a set of predefined categories listed in Fig. 4. Section 12 contains detailed descriptions of each category, the prompts used to implement the categorization with 1.5 Flash, and the correlation between 1.5 Flash category labels and those of SOTA models.

| Category | Description |
|---|---|
| *Multi-Turn* | Conversations with multiple turns. |
| *Exclude Ties* | Battles which do not end in a tie. |
| *Exclude Refusal* | Neither model refuses to answer. |
| *Captioning* | *Only* asks for a description of the image. |
| *OCR* | Requires reading text within the image. |
| *Coding* | Contains a code block in either the user inputs or model outputs. |
| *Entity Recognition* | Asks to identify objects, places, or people in the image. |
| *Homework* | Requires answering a problem which likely comes from a homework or exam. |
| *Humor* | Asks to explain the humor within the image or ask for a humorous composition. |
| *Diagram* | Contains images with a diagram (e.g., flowchart, circuit, graph). |
| *Creative Writing* | Asks for a creative composition such as a story or a script. |

Figure 4. Descriptions of VisionArena categories.



Figure 5. **Category Distribution.** Excluding preset examples. We see that direct chat data contains a higher proportion of coding, homework, and diagram questions while battle data contains more captioning, humor, and creative writing questions.



Figure 6. **Impact of confounding variables on user preferences, measured by $\hat{\gamma}$ in the enhanced Bradley-Terry Model.** Length is by far the most influential stylistic factor, with higher influence on preference for more open ended questions like humor, creative writing, and captioning.

Fig. 5 shows the distribution of category counts for both battles and direct chat conversations. We observe that direct chat conversations contain more homework and diagram understanding problems while battles contain more humor, captioning, and creative writing problems. We speculate this is because users in direct chat mode are more interested in using proprietary VLMs at no cost to assist them with their daily tasks.

Model rankings and Arena Scores for select categories and languages are shown in Fig. 7 (full table and arena scores can be found in the supplemental). We find several interesting insights such as:

- Gemini 1.5 Pro Exp [37], InternVL [7, 8] and Reka Flash [42] achieve a worse ranking for categories with require OCR like coding, homework, and diagrams.
- InternVL has a large improvement in ranking for the captioning category. In the following section, we show that this is largely due to stylistic choices such as formatting.
- Claude Opus, Sonnet, and Haiku [3] see an increase in ratings for multi-turn conversations.
- The Gemini class of models drops in performance on non-english conversations.

## 4.4. Controlling stylistic biases in evaluations

The Vision Arena captures signals from users of various backgrounds and preferences to construct its leaderboard. However, recent literature has pointed out potential confounding variables in model evaluation such as the length of the response or stylistic formatting [9, 11]. Others have also mentioned various axes in which annotators may disagree including task underspecification, response style, refusals, and annotation errors [46]. Thus, we explore the effect of these stylistic features on the VisionArena user preference.

We follow recent work that extends the BT model to include style features [23]. Given a set of style features (e.g. response length, number of markdown headers), we add a style vector to the BT model $\vec{Z}$ where= $Z_i \in \mathbb{R}^S$ is a vector of $S$ style features. the enhanced BT model has the style coefficients $\gamma \in \mathbb{R}^S$:

$$\hat{\beta}, \hat{\gamma} = \arg \min_{\beta \in \mathbb{R}^M, \gamma \in \mathbb{R}^S} \frac{1}{n} \sum_{i=1}^{n} \text{CE}(\sigma(X_i^\top \beta + Z_i^\top \gamma), Y_i)$$

For each style feature $Z_i$, we compute the normalized difference between the feature values of both model responses. The resulting $\hat{\beta}$ represents model strengths adjusted for style effects, while $\hat{\gamma}$ quantifies the influence of style on user preferences.

To control for these stylistic factors, we modify how VisionArena computes the model scores by accounting for the stylistic differences between two answers (response length, number of markdown headers, etc) as additional features to the existing BT model.

**Figure 7. Model rankings across question categories and languages.** Cells with * have fewer than 100 votes. Certain models achieve a much higher ranking for a particular category, such as InternVL2 on captioning and Chinese, Reka Core on entity recognition, and Llava 1.6 on Humor. Conversely, we see certain model rankings drop, such as Reka Flash on multi-turn and diagrams, and Gemini on Chinese.

| | Overall | Multi-Turn | Exclude Refusal | Captioning | OCR | Entity Recognition | Coding | Homework | Diagram | Humor | Creative Writing | English | Chinese | Russian | Vietnamese |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Gemini 1.5 Pro Exp | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 1 | 2 | 1 | 1 | 1 | 3 | 2 | 2 |
| GPT-4o | 2 | 2 | 2 | 3 | 1 | 1 | 1 | 2 | 1 | 2 | 2 | 2 | 1 | 3 | 1 |
| Claude 3.5 Sonnet | 3 | 3 | 3 | 2 | 3 | 3 | 3 | 4 | 3 | 3 | 4 | 3 | 2 | 1 | 3 |
| Gemini 1.5 Pro | 4 | 4 | 4 | 4 | 4 | 4 | 5 | 3 | 5 | 5 | 5 | 4 | 7 | 4 | 4 |
| GPT-4 Turbo | 5 | 5 | 5 | 5 | 5 | 5 | 4 | 5 | 4 | 4 | 3 | 5 | 6 | 5 | 6 |
| GPT-4o Mini | 6 | 6 | 6 | 6 | 6 | 7 | 6 | 6 | 6 | 6 | 6 | 6 | 5 | 6 | 5 |
| Gemini 1.5 Flash | 7 | 9 | 7 | 8 | 7 | 6 | 7 | 7 | 7 | 8 | 7 | 7 | 9 | 8 | 8 |
| Claude 3 Opus | 8 | 7 | 8 | 9 | 8 | 8 | 8 | 9 | 8 | 7 | 8 | 9 | 8 | 7 | 9 |
| InternVL2 26b | 9 | 10 | 10 | 7 | 10 | 12* | 9 | 11 | 10 | 10 | 10 | 8 | 4 | 14 | 13 |
| Claude 3 Sonnet | 10 | 8 | 9 | 10 | 9 | 11 | 10 | 8 | 9 | 12 | 9 | 10 | 10 | 9 | 7 |
| Reka Flash Preview | 11 | 14 | 12 | 11 | 13 | 10 | 14 | 14 | 15 | 14 | 11 | 15 | 13 | 10 | 11 |
| Claude 3 Haiku | 12 | 11 | 11 | 15 | 11 | 13 | 11 | 10 | 11 | 16 | 12 | 16 | 11 | 11 | 10 |
| Reka Core | 13 | 13 | 13 | 14 | 12 | 9 | 12 | 12 | 12 | 13 | 13 | 13 | 14 | 12 | 12 |
| LLAVA 1.6 34b | 14 | 12 | 14 | 13 | 14 | 16 | 13 | 13 | 13 | 9 | 16 | 11 | 15 | 16 | 16* |
| MiniCPM v2.6 | 15 | 15 | 15 | 16 | 15 | 14* | 15* | 15 | 14 | 15* | 14* | 12 | 12* | 15* | 17* |
| CogVLM2 Llama3 Chat 19b | 16 | 16 | 16 | 12 | 16 | 15* | 17* | 16 | 16 | 11* | 15* | 14 | 16* | 13* | 14* |
| Phi 3 Vision 128k Instruct | 17 | 17* | 17 | 17 | 17 | 17* | 16* | 17* | 17 | 17* | 17* | 17 | 17* | 17* | 15* |



**Figure 8. Change in model rankings when style control is applied.** Cells with * contain battles with fewer than 100 votes. InternVL, Reka Flash, and Llava 1.6 see a ranking drop across the majority of categories, while Claude 3 Haiku and MiniCPM see an increase across most categories. The model lineup changes the most for the captioning, humor, and Vietnamese categories.

| | Overall | Multi-Turn | Exclude Ties | Exclude Refusal | Captioning | OCR | Coding | Entity Recognition | Homework | Humor | Diagram | Creative Writing | English | Chinese | Russian | Vietnamese |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Gemini 1.5 Pro Exp | -1 | 0 | -1 | -1 | -2 | 0 | 0 | 0 | -1 | 0 | -1 | 0 | -1 | 0 | 0 | -1 |
| GPT-4o | 1 | 0 | 1 | 1 | 2 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | -1 |
| Claude 3.5 Sonnet | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | -1 | 0 | 0 | 0 | 2 |
| Gemini 1.5 Pro | 0 | 0 | 0 | 0 | -1 | 0 | 1 | 0 | 0 | 1 | 1 | 2 | 0 | 1 | 0 | 0 |
| GPT-4 Turbo | 0 | 0 | 0 | -1 | -1 | 0 | -1 | 0 | 0 | -1 | -1 | -1 | 0 | -1 | -1 | -3 |
| GPT-4o Mini | 0 | 0 | 0 | 1 | 2 | 0 | 0 | 0 | -1 | 0 | 0 | 0 | 1 | 1 | 0 | -1 |
| Gemini 1.5 Flash | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 |
| Claude 3 Opus | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | -1 | 0 | 0 | 1 | -1 | 0 | 1 |
| InternVL2 26b | -2 | -1 | -1 | 0 | -3 | -1 | -1 | -2* | -1 | -3 | -3 | 0 | -1 | -1 | -1 | -1 |
| Claude 3 Sonnet | 1 | 0 | 1 | 0 | 1 | 0 | 1 | -1 | -1 | 3 | 0 | 0 | 0 | -1 | 0 | 2 |
| Reka Flash Preview | -3 | -1 | -3 | -2 | -3 | -2 | 0 | -1 | -1 | -2 | -1 | -1 | -1 | -1 | -3 | -1 |
| Claude 3 Haiku | 2 | 1 | 2 | 0 | 3 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 4 | 1 | 1 | 0 |
| Reka Core | 1 | 1 | 0 | 1 | 3 | 0 | 0 | 0 | 1 | 2 | 0 | -1 | 0 | 1 | 1 | 1 |
| LLAVA 1.6 34b | -1 | -2 | -1 | -1 | -2 | 0 | 0 | 0 | -1 | -3 | -1 | 0 | -3 | 0 | -1 | -1* |
| MiniCPM v2.6 | 2 | 2 | 1 | 2 | 0 | 2 | 0* | 4* | 2 | 1* | 3 | 1* | 1 | 0* | 1* | 1* |
| CogVLM2 Llama3 Chat 19b | 0 | 0 | 1 | 0 | -1 | 0 | 0* | 0* | 0 | 1* | -1 | 0* | -1 | 0* | 1* | -1* |
| Phi 3 Vision 128k Instruct | 0 | 0* | 0 | 0 | 0 | 0 | 0* | 0* | 0* | 0* | 2 | 0* | 0 | 0* | 1* | 2* |

**Controlling for length and markdown.** Applying style control, we see the captioning category is heavily affected by style, with Fig. 8 showing a large difference in model rankings. We suspect this is because VLMs are heavily optimized for captioning and can usually correctly identify the main subjects and context of the image. This is supported in Table 2, which shows captioning questions have the smallest proportion of 'both bad' votes. In cases where both models provide a reasonable description of the image, the user may rely on stylistic features to determine preference.
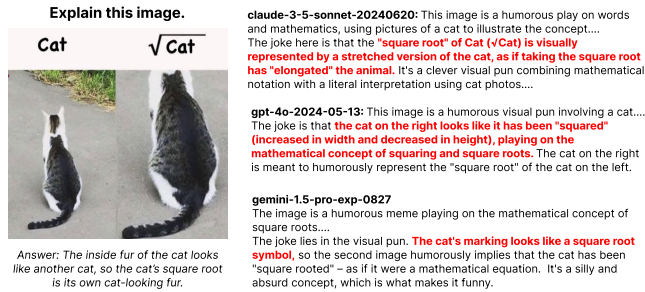
Furthermore, models like InternVL and Reka Flash Preview see a large decrease in rankings when style control is applied. Fig. 13 shows these models have an unusually high output token count compared to models of a similar class, indicating that some models may be 'preference hacking' by training their models to produce long or nicely-formatted outputs. While this is not necessarily bad, it is important to

consider when decoupling preference from capability.

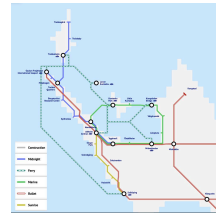| Caption | Homework | OCR | Coding | Humor | Entity Recog. | Creative Writing | Diagram |
|---|---|---|---|---|---|---|---|
| 11.33 | 31.12 | 20.62 | 29.4 | 18.47 | 22.25 | 14.21 | 26.06 |

Table 2. **Percentage of 'tie (bothbad)' per category.** Captioning, creative writing, and humor categories have low percentages of bothbad responses compared to coding and homework problems.

**Controlling for specificity.** We further extend the Bradley-Terry model to include the effect of response specificity. We define the complexity of a response as the number of named entities in the response. We use a NER model [32] to tag each response and use the number of named entities as our specificity score. In Fig. 6, we see that users prefer high specificity for tasks like entity recognition and diagram understanding while placing less emphasis on specificity for tasks like captioning and homework.

**Explain this image.**

**claude-3-5-sonnet-20240620:** This image is a humorous play on words and mathematics, using pictures of a cat to illustrate the concept.... The joke here is that the **"square root" of Cat (√Cat) is visually represented by a stretched version of the cat, as if taking the square root has "elongated" the animal.** It's a clever visual pun combining mathematical notation with a literal interpretation using cat photos....

**gpt-4o-2024-05-13:** This image is a humorous visual pun involving a cat.... The joke is that **the cat on the right looks like it has been "squared" (increased in width and decreased in height), playing on the mathematical concept of squaring and square roots.** The cat on the right is meant to humorously represent the "square root" of the cat on the left.

**gemini-1.5-pro-exp-0827**
The image is a humorous meme playing on the mathematical concept of square roots....
The joke lies in the visual pun. **The cat's marking looks like a square root symbol,** so the second image humorously implies that the cat has been "square rooted" – as if it were a mathematical equation. It's a silly and absurd concept, which is what makes it funny.

*Answer: The inside fur of the cat looks like another cat, so the cat's square root is its own cat-looking fur.*

(a) **Failure case 1.**

**Which line connects Trolskaholm to Solköping?**

**claude-3-5-sonnet-20240620:** The line that connects Trolskaholm to Solköping is the **Bullet line,** which is represented by the red line on the map. **This red line runs from Trolskaholm through several stations including Alfheim Central and Singer Park,** before reaching Solköping in the southern part of the map.

**gpt-4o-2024-05-13:** The line that connects **Trolskaholm to Solköping is the "Midnight" line,** which is represented by the blue color on the map.

**gemini-1.5-pro-exp-0827:** The **\*\*Sunrise\*\* line** connects Trolskaholm to Solköping. It's the yellow/gold line on the map. **You can follow it from Trolskaholm all the way down to Solköping.**

*Answer: The Ferry Line*

(b) **Failure case 2.**

Figure 9. **VLM failure modes.** The top proprietary models fail on questions which require advanced visual reasoning. For example, failure case 1 requires the visual understanding that the cat's fur patterns looks like another smaller black cat, and the linguistic connection between this and a square root. Non-truncated outputs in Section 16.

## 4.5. Failure Cases

We use VisionArena to analyze examples which are particularly challenging for current VLMs. We first filter VisionArena-Chat for prompts where the user voted that both models are bad, and collect a pool of 10 images which most or all of the current VLMs fail. The full set can be found in the supplemental section. Fig. 9 shows two examples of user questions that require advanced visual reasoning. This example requires the model to (1) understand that the two cats in the image are the same (2) the pattern on the cats back gives the illusion of another cat (3) this illusion of a smaller cat in a cat can be related to the fact that the square root of a number is a smaller demonimation of that number. Fig. 9b is an example of fine-grained spatial understanding, as the model must locate both locations and reason over the many intersecting lines in the image. Analysis of other failure cases in the supplemental section indicates that current VLMs still struggle on visual grounding tasks like reading distorted images, spatial understand and counting, as well as more complex reasoning tasks.

## 5. Instruction tuning vision-language models

Effective instruction finetuning for vision-language models depends on the diversity of instructions, the difficulty of prompts, and the quality of responses. This section demonstrates the potential of VisionArena for training high-performance instruction-following models.

We curate a high quality instruction-tuning dataset by sampling from the conversations with the highest-performing VLMs. We choose 100,000 conversations from VisionArena-Chat from the top models. This led to a dataset with conversations from 16 different models[1] including proprietary models such as GPT-4o [33], Gemini-1.5-Pro [37],

and Claude-3.5-Sonnet [3] as well as open-source models such as Qwen2-VL-72B [5, 43] and Llama-3.2-90B-Vision-Instruct [2]. We compare the effectiveness of this dataset for finetuning to a 100K subset of the Llava-Instruct-158K [24].

We use Llama-3.2-11B-Vision and freeze the vision encoder while finetuning the multimodal projector and language model. We finetune for 3 epochs on the data for both our 100k dataset and the 100k Llava-Instruct dataset. In Table 3, we label the model finetuned on our VisionArena data Llama-3.2-VisionArena and the model trained using Llava-Instruct-158K Llama-3.2-Llava-Instruct.

The evaluation results are shown in Table 3. Llama-3.2-VisionArena significantly outperforms Llama-3.2-Llava-Instruct on MME, HallusionBench, MMMU, MMMU-Pro, and WildVision-Bench. Additionally, Llama-3.2-VisionArena outperforms Llama-3.2-11B-Vision-Instruct on both MME and WV-Bench, despite being fine-tuned on 30X less data. In Sec. 13, we show that these improvements are not due to contamination of these benchmarks

## 6. VisionArena-Bench: An Automatic Offline Human-Preference Benchmark for VLMs

Lastly, we demonstrate VisionArena's ability to cheaply approximate model preference rankings with VisionArena-Bench. Currently, online preference benchmarks like Chatbot Arena obtain a ranking for a new model by adding it to their platform and waiting days or weeks to collect enough votes for a stable ranking. For a single developer hoping to test a particular version of their model, obtaining these online model rankings is infeasible.

We develop a solution for those who need a quick and cheap evaluation of their models: VisionArena-Bench. VisionArena-Bench is a set of 500 diverse image and text

---

[1] gpt-4o-mini-2024-07-18, gpt-4-turbo-2024-04-09, gemini-1.5-pro-api-0514, claude-3-5-sonnet-20240620, gemini-1.5-pro-exp-0827, gpt-4o-2024-05-13, gemini-1.5-pro-exp-0801, gemini-1.5-flash-api-0514, claude-3-opus-20240229, gemini-1.5-flash-exp-0827, gemini-1.5-flash-8b-exp-0827, llama-3.2-vision-90b-instruct, qwen2-vl-72b [43], gpt-4o-2024-08-06, chatgpt-4o-latest-20240903, chatgpt-4o-latest-20240808

[2] gemini-1.5-pro-exp-0827, gemini-1.5-flash-exp-0827[37], gpt-4o-2024-05-13, gemini-1.5-flash-8b-exp-0827, internvl2-26b [8], claude-3.5-sonnet-20240620 [3], gpt-4-turbo-2024-04-09 [33], claude-3-sonnet-20240229, llama-3.2-11b-vision-instruct [2], gemini-1.5-pro-001, internvl2-4b, gpt-4o-mini-2024-07-18, claude-3-opus-20240229, gemini-1.5-flash-001, reka-core-20240501 [42], claude-3-haiku-20240307

| | | MME[12] | | HallucinationBench[17] | | | MMMU[44] | MMMU-Pro[45] | WV-Bench[26] |
|---|---|---|---|---|---|---|---|---|---|
| Model | # Samples | Cog. | Perc. | Acc.(all) | Fig. | Q. | Acc. | Acc. | Acc. |
| Llama3.2-11B-V-Instruct | 3M+ | <u>327.5</u> | <u>1421.7</u> | **48.6** | **26.0** | **23.1** | **50.7** | **0.28** | <u>47.2</u> |
| Llama-3.2-Llava-Instruct | 100K | 262.1 | 1067.6 | 38.7 | 13.9 | 8.6 | 27.9 | 0.12 | 10.4 |
| Llama-3.2-VisionArena | 100K | **345.4** | **1437.0** | <u>45.2</u> | <u>19.9</u> | <u>16.3</u> | <u>43.0</u> | <u>0.27</u> | **56.9** |

Table 3. **Performance across models trained with different instruction tuning datasets.** Llama-3.2-11B-Vision-Instruct scores are author-reported. Fine-tuning on samples from VisionArena-Chat outperforms fine-tuning on Llava-Instruct across all benchmarks. Llama-3.2-VisionArena also outperforms Llama-3.2-11B-Vision-Instruct on both MME and WV-Bench, despite being fine-tuned on 30x less data.

prompts that accurately approximates the model ranking from the Chatbot Arena VLM Leaderboard.

*Offline benchmark curation.* To gather questions for this offline benchmark, we build on recent work in building benchmarks from crowd-sourced evaluations in LLM and adapt them to the context of VLM [23]. We first filter the data for single turn to prevent the user from correcting the model in its response on its second turn. We then filter out non-english conversations as we are personally unable to verify the quality of non-english prompts.

To sample diverse questions, we perform topic modeling using the library BERTopic with multimodal embeddings [16]. We extract image embeddings and text embeddings using a CLIP model (e.g. CLIP-ViT-B-32) [36]. We then average the image and text embeddings so that each document corresponds to a single embedding. Then, we use UMAP to reduce the dimensions of the embedding and use hierarchical-based clustering (HDBSCAN) to generate topic clusters [27, 30]. We then uniformly sample from each topic cluster to generate the 500 prompts.

*Automatic evaluation with VLM-as-a-judge.* To evaluate a model, we use the LLM-as-a-judge framework mentioned in [23] applied to VLMs. We first select a fixed anchor model (GPT-4-Turbo[33]) that will be used in the pairwise comparisons. To evaluate a given model $M$ on a user prompt $p$, we generate responses for both $M$ and the anchor model on $p$ and then utilize GPT-4o as a judge to provide a preference score between the $(M, \text{anchor})$ pair on a 5-point Likert scale. 1 indicates a strong preference for model A and 5 indicates a strong preference for model B. We then obtain this score for all models across all prompts in VisionArena-Bench to obtain VLM-generated pairwise preference votes and use the same procedure as described in Sec. 3 to produce final model scores. In Sec. 14, we provide the detailed judge prompt template. To avoid potential bias, we prompt the judge model to judge twice, swapping the response position between the two rounds.

To evaluate the effectiveness of our benchmark to existing work [26], we leverage standard metrics such as Spearman correlation and Kendall Tau correlation which measure the agreement between two benchmarks' model rankings. We choose a shared set of 16 models and compare the offline benchmark rankings to the online Chatbot Arena. Table 4 shows VisionArena-Bench achieves a higher Spearman and Kendall Tau correlation than WildVision-Bench, with a 17.1% and 20.5% gain respectively. In Sec-

tion 14, we also compare the results using the same baseline model as WildVision (e.g. Claude-3-Sonnet-20240229), showing that the spearman correlation to Chatbot Arena's VLM leaderboard (10/23/2024) remains the same. This demonstrates the potential of VisionArena-Bench as a cost-effective and scalable offline benchmark that closely mirrors human preferences captured in online evaluations, enabling researchers to efficiently assess and compare VLMs without the need for extensive user studies.

| | VisionArena-Bench | WV-Bench |
|---|---|---|
| Confidence Agreement | **98.6%** | 87.6% |
| Spearman Correlation | **97.3%** | 80.2% |
| Kendall Tau Correlation | **89.7%** | 69.2% |

Table 4. **Correlation of rankings with ChatbotArena's VLM leaderboard.** Performance comparison on 16 models[2] between VisionArena-Bench to Chatbot Arena's VLM leaderboard based on confidence agreement, spearman correlation, and Kendall tau correlation. Result as of leaderboard on October 23, 2024.

# 7. Discussion, Limitations, and Future Work

Human preference benchmarks provide a critical lens for assessing performance on open-ended tasks where an explicit notion of correctness is either unavailable or subjective. Users may implicitly consider factual accuracy when making their preferences, but we would like to emphasize that **this benchmark is designed to measure human preferences rather than explicitly evaluate factual accuracy.** We see VisionArena as complementary to existing datasets and benchmarks that measure objective correctness.

Despite the breadth of coverage offered by VisionArena, significant gaps remain in representing the full distribution of real-world use cases for vision-language models. As highlighted in Sec. 4, our dataset contains many sample from domains such as STEM problems, OCR tasks, and toy problems (e.g., humor and riddles). These areas, while valuable, leave critical application domains underrepresented, including geospatial applications, medical domains, and visual assistance. Furthermore, while VisionArena contains over 100 languages, many of these languages do not contain enough examples to produce a stable leaderboard. Looking forward, we hope to enourage a more diverse user base by changing our UI to be multi-lingual and improving general user experience. Lastly, we have made it easy for the community to contribute new question categories and models at https://github.com/lm-sys/FastChat.

# References

[1] Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou. Gradio: Hassle-free sharing and testing of ml models in the wild. *arXiv preprint arXiv:1906.02569*, 2019. 1

[2] AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md. 1, 7

[3] Anthropic. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf, 2024. (Accessed on 06/05/2024). 1, 5, 7

[4] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015. 1

[5] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023. 7

[6] Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. doi: 10.2307/2334029. URL https://doi.org/10.2307/2334029. 3

[7] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023. 5

[8] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 5, 7

[9] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: An open platform for evaluating llms by human preference, 2024. 1, 2, 3, 5

[10] Yuntian Deng, Wenting Zhao, Jack Hessel, Xiang Ren, Claire Cardie, and Yejin Choi. Wildvis: Open source visualizer for million-scale chat logs in the wild, 2024. URL https://arxiv.org/abs/2409.03753. 2

[11] Lisa Dunlap, Krishna Mandal, Trevor Darrell, Jacob Steinhardt, and Joseph E Gonzalez. Vibecheck: Discover and quantify qualitative differences in large language models. *arXiv preprint arXiv:2312.02974*, 2024. URL https://arxiv.org/abs/2410.12851. 5

[12] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models, 2024. URL https://arxiv.org/abs/2306.13394. 2, 8

[13] Google. Vision ai — cloud vision api. https://cloud.google.com/vision. Accessed: [date]. 3

[14] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1

[15] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering, 2017. URL https://arxiv.org/abs/1612.00837. 2

[16] Maarten Grootendorst. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*, 2022. 4, 8

[17] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 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 (CVPR)*, pages 14375–14385, June 2024. 8

[18] Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1

[19] Danna Gurari, Qing Li, Chi Lin, Yinan Zhao, Anhong Guo, Abigale J. Stangl, and Jeffrey P. Bigham. Vizwiz-priv: A dataset for recognizing the presence and purpose of private visual information in images taken by blind people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1

[20] Lalit Jain, Kevin Jamieson, Robert Mankoff, Robert Nowak, and Scott Sievert. The New Yorker cartoon caption contest dataset, 2020. URL https://nextml.github.io/caption-contest-data/. 3

[21] Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations – democratizing large language model alignment. *arXiv preprint arXiv:2304.07327*, 2023. URL https://doi.org/10.48550/arXiv.2304.07327. Published in NeurIPS 2023 Datasets and Benchmarks. 2

[22] Seongyun Lee, Seungone Kim, Sue Hyun Park, Geewook Kim, and Minjoon Seo. Prometheus-vision: Vision-language model as a judge for fine-grained evaluation, 2024. URL https://arxiv.org/abs/2401.06591. 3

[23] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks:

Arena-hard and benchbuilder pipeline, 2024. URL https://arxiv.org/abs/2406.11939. 5, 8

[24] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. URL https://arxiv.org/abs/2304.08485. 2, 7

[25] Yujie* Lu, Dongfu* Jiang, Hui* Chen, Xingyu Fu, Yingzi Ma, Jing Gu, Michael Saxon, Chaowei Xiao, Wenhu Chen, Yejin Choi, Bill Yuchen Lin, Miguel Eckstein, and William Wang. Wildvision data and model, 2024. URL https://huggingface.co/WildVision. 1

[26] Yujie Lu, Dongfu Jiang, Wenhu Chen, William Yang Wang, Yejin Choi, and Bill Yuchen Lin. Wildvision: Evaluating vision-language models in the wild with human preferences. 2024. 1, 2, 3, 8

[27] Claudia Malzer and Marcus Baum. A hybrid approach to hierarchical density-based cluster selection. In *2020 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI)*, page 223–228. IEEE, September 2020. doi: 10.1109/mfi49285.2020.9235263. URL http://dx.doi.org/10.1109/MFI49285.2020.9235263. 8

[28] Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.177. URL https://aclanthology.org/2022.findings-acl.177. 3

[29] Minesh Mathew, Dimosthenis Karatzas, R Manmatha, and CV Jawahar. Docvqa: A dataset for vqa on document images. corr abs/2007.00398 (2020). *arXiv preprint arXiv:2007.00398*, 2020. 2, 3

[30] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction, 2020. URL https://arxiv.org/abs/1802.03426. 8

[31] Microsoft. Photodna service. Online, 2024. URL https://www.microsoft.com/en-us/photodna. Accessed: 2024-11-18. 3

[32] Joel Nothman, Nicky Ringland, Will Radford, Tara Murphy, and James R. Curran. Learning multilingual named entity recognition from wikipedia. *Artificial Intelligence*, 194:151–175, 2013. ISSN 0004-3702. doi: https://doi.org/10.1016/j.artint.2012.03.006. URL https://www.sciencedirect.com/science/article/pii/S0004370212000276. Artificial Intelligence, Wikipedia and Semi-Structured Resources. 6

[33] OpenAI. Gpt-4 technical report, 2023. 1, 7, 8

[34] OpenAI. Openai text moderation tool. Online, 2024. URL https://platform.openai.com/docs/guides/moderation. Accessed: 2024-11-18. 3

[35] Karl Pearson. Note on regression and inheritance in the case of two parents. *Proceedings of the Royal Society of London*, 58:240–242, 1895. ISSN 03701662. URL http://www.jstor.org/stable/115794. 4

[36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,

Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL https://arxiv.org/abs/2103.00020. 8

[37] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 1, 5, 7

[38] Nimrod Shabtay, Felipe Maia Polo, Sivan Doveh, Wei Lin, M. Jehanzeb Mirza, Leshem Chosen, Mikhail Yurochkin, Yuekai Sun, Assaf Arbelle, Leonid Karlinsky, and Raja Giryes. Livexiv – a multi-modal live benchmark based on arxiv papers content, 2024. URL https://arxiv.org/abs/2410.10783. 3

[39] Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8317–8326, 2019. 3

[40] Charles Spearman. The proof and measurement of association between two things. *The American Journal of Psychology*, 15(1):72–101, January 1904. doi: 10.2307/1412159. URL https://doi.org/10.2307/1412159. 4

[41] Wei Ren Tan, Chee Seng Chan, Hernan Aguirre, and Kiyoshi Tanaka. Improved artgan for conditional synthesis of natural image and artwork. *IEEE Transactions on Image Processing*, 28(1):394–409, 2019. doi: 10.1109/TIP.2018.2866698. URL https://doi.org/10.1109/TIP.2018.2866698. 3

[42] Reka Team, Aitor Ormazabal, Che Zheng, Cyprien de Masson d'Autume, Dani Yogatama, Deyu Fu, Donovan Ong, Eric Chen, Eugenie Lamprecht, Hai Pham, Isaac Ong, Kaloyan Aleksiev, Lei Li, Matthew Henderson, Max Bain, Mikel Artetxe, Nishant Relan, Piotr Padlewski, Qi Liu, Ren Chen, Samuel Phua, Yazheng Yang, Yi Tay, Yuqi Wang, Zhongkai Zhu, and Zhihui Xie. Reka core, flash, and edge: A series of powerful multimodal language models, 2024. URL https://arxiv.org/abs/2404.12387. 5, 7

[43] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution, 2024. URL https://arxiv.org/abs/2409.12191. 7

[44] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi, 2024. URL https://arxiv.org/abs/2311.16502. 1, 2, 8

[45] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhu Chen, and Graham Neubig.

Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. *arXiv preprint arXiv:2409.02813*, 2024. 8

[46] Michael JQ Zhang, Zhilin Wang, Jena D. Hwang, Yi Dong, Olivier Delalleau, Yejin Choi, Eunsol Choi, Xiang Ren, and Valentina Pyatkin. Diverging preferences: When do annotators disagree and do models know?, 2024. URL https://arxiv.org/abs/2410.14632. 5

[47] Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Yin and Yang: Balancing and answering binary visual questions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1

[48] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric. P Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. Lmsys-chat-1m: A large-scale real-world llm conversation dataset, 2023. 2

# VisionArena: 230K Real World User-VLM Conversations with Preference Labels

## Supplementary Material

## 8. Acknowledgments

## 9. Interface Details

We implement our interface in Gradio [1]. If a user uploads an image in the first turn, two random VLMs are selected to answer the query. A user can only chat with one image per conversation. As shown in Fig. 10, a user can also select a random image which will select from our preset examples. Note that the user still needs to come up with a query, even for these preset images.

## 10. Topic Distribution

In Fig. 11 we show the top 20 topic clusters from VisionArena-Chat and WildVision-Chat. The topic clusters are extracted from 50K sampled English conversations from VisionArena-Chat and 37K English conversations (all) from WildVision-Chat. We see that VisionArena-Chat includes more diverse and broad topics especially in the STEM field. Furthermore, WildVision has very specific clusters like "detailed cartoon dog description", "Shock absorber assembly analysis", "rice lead disease identification". Looking at these clusters we see that they contain a large number of duplicate prompts.

## 11. More Data Stats

We provide further information on language distribution (Fig. 15, Fig. 16), battle counts (Fig. 12, Fig. 14), token count (Fig. 13), turn distribution (Fig. 18, Fig. 19), proportion of refusals (Fig. 20), battle outcome counts (Fig. 17), win rates (Fig. 21), and category overlap (Fig. 24).



Figure 12. **VisionArena-Battle counts per model.**



Figure 13. **Model token count in VisionArena-Battle.** Models in bold see a large decrease in rank when style control is applied.

Figure 10. **Interface for anonymous side-by-side chat.**



Figure 11. **Top 20 topic clusters of VisionArena-Chat compared to WildVision-Chat.** VisionArena-Chat includes more diverse and broad topics especially in the STEM field.

Figure 14. **VisionArena-Chat counts per model.**



Figure 15. **VisionArena-Battle counts for the top 25 languages.**



Figure 16. **VisionArena-Chat counts for the top 25 languages.**



Figure 17. **Battle Outcome Counts.**

Figure 18. **VisionArena-Battle Conversation Turn Distribution**



Figure 19. **VisionArena-Chat Conversation Turn Distribution**



Figure 20. **Proportion of Refusals per model.**

## 12. Category Details

Below are the system prompts used to classify user prompts into the categories described in Section 4. We find classifications are more accurate for certain categories when using only the prompt or only the image, as indicated in the prompt titles. We use Gemini 1.5 Flash for classification and show in Table 5 that our classifications achieve high agreement to using SOTA models as category labelers.

---

**OCR System Prompt (image + prompt)**

You are tasked with determining if a given VQA question is an optical character recognition (OCR) question. An OCR question requires reading and understanding text in the image to answer. If there is some amount of text in the image and the question requires reading the text in any capacity it should be classified as Optical Character Recognition.

Output your verdict in the following format: "<decision>[yes/no]<decision>". Do NOT explain.

---

**Refusal System Prompt (responses only)**

You are tasked with determining if any of the given model responses are a refusal to answer. A refusal occurs when the model explicitly declines to answer or complete a task, due to reasons like safety, moderation, or model limitations (e.g. the user asks the model to search the web but it does not have that capability). You will be given the responses of 2 models, A and B and you are to determine if A refuses to answer, B refuses, both refuse, or neither refuse.

Output your verdict in the following format: "<decision>[yes/no]<decision>". Do NOT explain.

---

**Captioning System Prompt (prompt only)**

You are tasked with determining if a given VQA question is a captioning question. A captioning question asks for a general, overall description of the entire image. It must be a single, open-ended query that does NOT ask about particular objects, people, or parts of the image, nor require interpretation beyond a broad description of what is visually present. Examples include 'What is happening in this image?', 'Describe this picture.', 'Explain', etc. An example of a non-captioning question is 'Describe what is funny in this picture.' because it asks for a specific interpretation of the image content.

Output your verdict in the following format: <decision>[yes/no]<decision>. Do NOT explain.

---

**Homework System Prompt (image only)**

You are tasked with determining if the given image contains a homework or exam question. A homework or exam question typically contains text with a well-defined question or task which asks for a solution. In addition, many homework and exam questions contain multiple choice, equations, and question numbers. You may also see text referring to showing your work or providing justification. Note that documents such as resumes, business cards, records, or personal notes are NOT considered homework or exam questions; homework and exam questions explicitly ask for a solution or explanation.
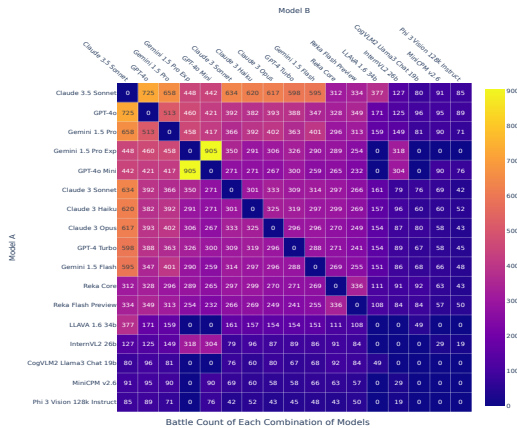
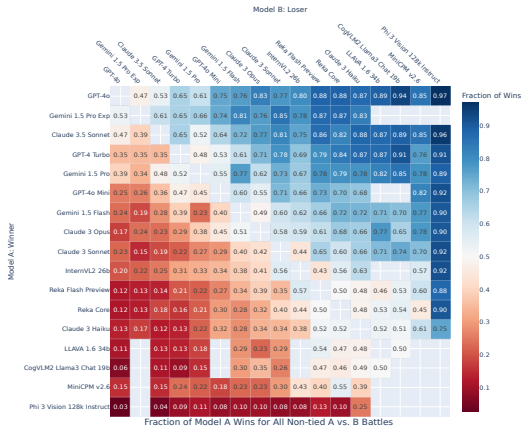Output your verdict in the following format: "<decision>[yes/no]<decision>". Do NOT explain.

(a) Bootstrap ELO Ratings



(b) Average Win Rate



(c) Battle Counts



(d) Win Fractions

Figure 21. **VisionArena-Battle Model Ranking Results.**

| Category | Uses Image | Labeler | Accuracy | Precision | Recall |
|----------|------------|---------|----------|-----------|--------|
| Homework | Yes | gemini-1.5-pro-exp-0827 | 0.987 | 0.929 | 0.967 |
| Captioning | No | claude-3-5-sonnet-20240620 | 0.967 | 0.938 | 0.934 |
| Humor | Yes | gemini-1.5-pro-exp-0827 | 0.925 | 0.788 | 0.636 |
| OCR | Yes | gemini-1.5-pro-exp-0827 | 0.818 | 0.954 | 0.769 |
| Entity Recognition | No | claude-3-5-sonnet-20240620 | 0.952 | 0.728 | 0.830 |
| Creative Writing | No | claude-3-5-sonnet-20240620 | 0.964 | 0.680 | 0.810 |
| Diagram | Yes | gemini-1.5-pro-exp-0827 | 0.961 | 0.858 | 0.953 |

Table 5. **Comparing Gemini-1.5-Flash question categorization to larger models.** Gemini-1.5-Flash is evaluated against SOTA models on 1000 samples from VisionArena-Chat, using Gemini-1.5-Pro for image-based prompts and Claude-3.5-Sonnet for text-based prompts. Gemini-1.5-Flash achieves high agreement with SOTA models for category classification.

**Humor Systems Prompt (image + prompt)**

You are tasked with determining if a given VQA question is a humor question. A humor question asks for a humorous or funny response based on the image or asks to understand what is funny about an image. This includes questions that

ask to explain an image which is humorous, such as memes.

Output your verdict in the following format:
"<decision>[yes/no]<decision>". Do NOT explain.

Figure 22. **Random Samples from VisionArena-Battle**



Figure 23. **Random Samples from VisionArena-Chat**

Figure 24. **VisionArena-Battle category overlap.**

## 13. Contamination with Existing Benchmarks

To ensure that our results from Sec. 5 are not due to training on questions from the test sets, we investigate the rate of benchmark contamination in VisionArena-Chat. Using OpenAI's text-embedding-small embeddings, we compute the cosine similarity between each VisionArena-Chat question and all benchmark questions, selecting the nearest neighbor with the highest similarity score. We then count the number of cases where this similarity is $\geq 0.8$, indicating minor rephrasings of the same question. Table 6 shows that less than 2% of benchmark questions are seen on VisionArena-Chat.

| Dataset | # Matches | % dataset | % VisionArena-Chat |
|---|---|---|---|
| MMMU | 47 | 0.4% | 0.02% |
| MME | 0 | 0.0% | 0.0% |
| HallusionBench | 0 | 0.0% | 0.00% |
| MMMU Pro | 23 | 1.3% | 0.01% |

Table 6. **Proportion of benchmark data in VisionArena-Chat.**

## 14. VisionArena-Bench

| Model | Score | 95% CI | Token # |
|---|---|---|---|
| gpt-4o-2024-05-13 | 67.7 | (-1.7, 1.8) | 316 |
| gemini-1.5-pro-exp-0827 | 66.2 | (-1.8, 1.5) | 329 |
| gemini-1.5-flash-exp-0827 | 60.3 | (-1.9, 1.9) | 367 |
| claude-3.5-sonnet-20240620 | 54.5 | (-2.1, 1.9) | 262 |
| gpt-4-turbo-2024-04-09 | 50.0 | (0.0, 0.0) | 258 |
| gemini-1.5-pro-001 | 45.5 | (-1.8, 2.0) | 261 |
| gpt-4o-mini-2024-07-18 | 40.0 | (-2.3, 1.9) | 224 |
| gemini-1.5-flash-8b-exp-0827 | 30.6 | (-2.3, 1.8) | 354 |
| internvl2-26b | 23.3 | (-2.1, 1.1) | 515 |
| gemini-1.5-flash-001 | 23.0 | (-1.1, 1.6) | 271 |
| claude-3-opus-20240229 | 18.9 | (-1.9, 1.7) | 201 |
| claude-3-sonnet-20240229 | 18.4 | (-1.4, 1.3) | 205 |
| reka-core-20240501 | 15.6 | (-1.3, 1.4) | 252 |
| llama-3.2-11b-vision-instruct | 11.2 | (-1.3, 1.1) | 457 |
| claude-3-haiku-20240307 | 9.6 | (-1.1, 1.0) | 155 |
| internvl2-4b | 6.8 | (-0.9, 0.8) | 421 |

Table 7. VisionArena-Bench leaderboard (baseline: GPT-4-Turbo)

| Model | Score | 95% CI | Token # |
|---|---|---|---|
| gemini-1.5-pro-exp-0827 | 87.6 | (-1.1, 1.1) | 329 |
| gpt-4o-2024-05-13 | 86.8 | (-1.1, 1.3) | 316 |
| claude-3-5-sonnet-20240620 | 86.3 | (-1.3, 1.2) | 262 |
| gemini-1.5-flash-exp-0827 | 83.5 | (-1.7, 1.1) | 367 |
| gpt-4-turbo-2024-04-09 | 80.7 | (-1.0, 1.6) | 258 |
| gemini-1.5-pro-001 | 75.3 | (-1.7, 1.4) | 261 |
| gpt-4o-mini-2024-07-18 | 73.0 | (-1.3, 1.4) | 224 |
| gemini-1.5-flash-8b-exp-0827 | 64.7 | (-1.5, 2.4) | 354 |
| gemini-1.5-flash-001 | 58.4 | (-1.9, 1.5) | 271 |
| internvl2-26b | 54.2 | (-1.7, 1.7) | 515 |
| claude-3-opus-20240229 | 52.0 | (-2.0, 1.7) | 201 |
| claude-3-sonnet-20240229 | 50.0 | (0.0, 0.0) | 205 |
| reka-core-20240501 | 37.9 | (-1.9, 1.7) | 252 |
| llama-3.2-11b-vision-instruct | 32.6 | (-1.7, 1.8) | 457 |
| claude-3-haiku-20240307 | 30.7 | (-2.3, 1.6) | 155 |
| internvl2-4b | 19.6 | (-1.9, 1.3) | 421 |

Table 8. VisionArena-Bench leaderboard (baseline: claude-3-sonnet-20240229)

## 15. Additional model details

Table 9 shows the mapping from the model names used in Section 4 to the exact model versions.

| Model Version | Model Name |
|---|---|
| claude-3-5-sonnet-20240620 | Claude 3.5 Sonnet |
| claude-3-haiku-20240307 | Claude 3 Haiku |
| claude-3-opus-20240229 | Claude 3 Opus |
| claude-3-sonnet-20240229 | Claude 3 Sonnet |
| cogvlm2-llama3-chat-19b | CogVLM2 Llama3 Chat 19b |
| gemini-1.5-flash-api-0514 | Gemini 1.5 Flash |
| gemini-1.5-pro-api-0514 | Gemini 1.5 Pro |
| gemini-1.5-pro-exp-0801 | Gemini 1.5 Pro Exp |
| gpt-4-turbo-2024-04-09 | GPT-4 Turbo |
| gpt-4o-2024-05-13 | GPT-4o |
| gpt-4o-mini-2024-07-18 | GPT-4o Mini |
| internvl2-26b | InternVL2 26B |
| llava-v1.6-34b | LLAVA 1.6 34b |
| minicpm-v-2_6 | MiniCPM v2.6 |
| phi-3-vision-128k-instruct | Phi 3 Vision 128k Instruct |
| reka-core-20240501 | Reka Core |
| reka-flash-preview-20240611 | Reka Flash Preview |

Table 9. Model Name to exact model version

## 16. Failure Cases

*Hard OCR* (Fig. 25). While VLMs perform well at transcribing easily legible text, they struggle with perturbed text (e.g., rotations, blur). Reading such difficult text is essential for real-world applications. We show two failure cases: one with unclear handwriting and another with rotated text.
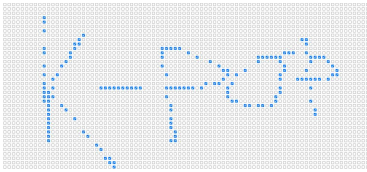
| Model | K-Pop | Sign | Shapes | Triangles | Meme | Map | Shoes | Chess | ARC |
|---|---|---|---|---|---|---|---|---|---|
| gemini-1.5-pro-exp-0827 | X | X | X | X | X | X | X | X | X |
| gpt-4o-2024-05-13 | X | O | X | X | X | X | O | X | X |
| claude-3-5-sonnet-20240620 | X | X | X | X | X | X | X | X | X |
| claude-3-opus-20240229 | X | X | X | X | X | X | O | X | X |
| gpt-4-turbo-2024-04-09 | X | X | X | X | X | X | X | X | X |
| gpt-4o-mini-2024-07-18 | X | X | X | X | X | X | X | X | X |
| gemini-1.5-pro-001 | X | X | X | X | X | X | X | X | X |
| gemini-1.5-flash-8b-exp-0827 | X | X | X | X | X | X | X | X | X |
| gemini-1.5-flash-exp-0827 | X | O | O | O | X | X | X | X | X |
| internvl2-26b | X | X | X | X | X | X | X | X | X |
| gemini-1.5-flash-001 | X | X | O | X | X | X | X | X | X |
| claude-3-sonnet-20240229 | X | X | X | X | X | X | O | X | X |
| llama-3.2-11b-vision-instruct | X | O | X | O | X | X | X | X | X |
| claude-3-haiku-20240307 | X | X | X | X | X | X | X | X | X |
| internvl2-4b | X | X | X | X | X | X | X | X | X |

Table 10. **Model performance across several hard tasks.** O indicates that the model solves the problem and X indicates that the model fails to solve the problem. 9 out of 16 models fail all questions.

*Counting* (Fig. 26). Counting is a critical skill for decision-making across education, organization, and daily life. While humans count effortlessly, VLMs still struggle. We provide two examples where the top three models fail: one requiring counting based on shape and color, and another involving intersecting triangles.

*Reasoning* (Fig. 27, Fig. 28). Reasoning is essential for helping users tackle complex problems. While reasoning remains a challenge for both VLMs and LLMs, we present five unique VLM failure cases.

**What does it put here?**



*Answer: K-Pop*

**gemini-1.5-pro-exp-0827:** It puts the word **"HELLO"** in a grid using blue checkmarks.

**gpt-4o-2024-05-13:** The image appears to be a **Minesweeper game board.** The blue numbers indicate the number of mines adjacent to each square. The numbers visible in the image are: - 1 -2 - 3 - 4 These numbers help the player deduce where the mines are located on the board. The goal of the game is to clear all the non-mine squares without detonating any mines.

**claude-3-5-sonnet-20240620:** The image shows a grid of small squares, with some squares filled in blue to form letters. The blue squares spell out **"THE LOTION"** in a pixelated or dot-matrix style font across the grid.

(a) K-Pop Failure Case

**What is written?**



*Answer: 6 Avenue West End*

**gemini-1.5-pro-exp-0827:** 9 West End Avenue

**gpt-4o-2024-05-13:** The text on the sign reads **"6 Avenue West End."**

**claude-3-5-sonnet-20240620:** The image shows text that reads **"West 6 Avenue"** when read from right to left. The text appears to be on a street sign or similar type of signage. The lettering is white on a dark blue background, and the image has a slightly blurred or distorted quality to it, possibly due to the angle or lighting conditions when the photo was taken.

(b) Sign Failure Case

Figure 25. **Hard OCR Failure Cases**. These failure cases show that models still fail in cases where the text is perturbed such as rotations or messy hand-writing.

**how many of each shape and color**



*Answer: 3 stars, 6 rectangles, 3 ovals*
*5 green shapes, 3 orange shapes, 4 blue shapes*

**gemini-1.5-pro-exp-0827:** Here's a count of the shapes and colors:
**Shapes:**
* **Stars:** 3
* **Ovals:** 3
* **Rectangles:** 5
**Colors:**
* **Light Blue:** 4 (3 stars, 1 rectangle)
* **Light Green:** 4 (2 ovals, 2 rectangles)
* **Light Orange/Peach:** 2 (1 oval, 1 rectangle)
Let me know if you have any other questions!

**gpt-4o-2024-05-13:** In the image, there are various shapes and colors. Here is the count for each shape and color:
**Shapes:**
- **Stars: 4**
- Ovals: 3
- Rectangles: 6
**Colors:**
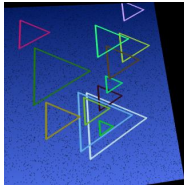- **Blue: 5 (4 stars, 1 rectangle)**
- Green: 5 (3 ovals, 2 rectangles)
- Orange: 3 (1 oval, 2 rectangles)

**claude-3-5-sonnet-20240620:** Here's a breakdown of the shapes and colors in the image:
Shapes:
- Stars: 3 (all light blue)
- Circles/Ovals: 3 (2 green, 1 peach/light orange)
- **Rectangles: 5** (3 green, 1 light blue, 1 peach/light orange)

Colors:
- Light Blue: 4 shapes (3 stars, 1 rectangle)
- Green: 5 shapes (2 circles, 3 rectangles)
- **Peach/Light Orange: 2 shapes (1 oval, 1 rectangle)**

In total, there are 11 shapes in the image.

(a) Shapes Failure Case

**How many distinct triangles there are on the image? Intersections between figures dont count.**



*Answer: 13*

**gemini-1.5-pro-exp-0827:** Let's count the distinct triangles in the image: 1. Top left: 1 (pink) 2. Top right: 1 (light purple) 3. Middle left: 1 (green) 4. Center: - Largest: 1 (light blue) - Medium: 1 (white) - Smallest: 1 (bright green) - Overlapping: 1 (brown) - Tiny: 1 (dark red) 5. Right: 1 (light blue)

Total count: 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 = 9

There are **9** distinct triangles in the image.

**gpt-4o-2024-05-13:** There are **11** distinct triangles in the image.

**claude-3-5-sonnet-20240620:** In the image, there are **10** distinct triangles. These triangles are of various sizes and colors, scattered across a blue background. The triangles appear in different colors including white, green, pink, yellow, and light blue. Some of the triangles are overlapping, but as per your instruction, we're not counting intersections between figures. Each triangle is considered as a separate, distinct shape regardless of its position or overlap with others.

(b) Triangles Failure Case

Figure 26. **Hard Counting Failure Cases**. These cases show that proprietary models still fail at counting tasks involving shapes even when the format is clear.

**Which line connects Trolskaholm to Solköping?**

**gemini-1.5-pro-exp-0827:** The **Sunrise** line connects Trolskaholm to Solköping.

It's the yellow/gold line on the map. **You can follow it from Trolskaholm all the way down to Solköping.**

**gpt-4o-2024-05-13:** The line that connects **Trolskaholm to Solköping is the "Midnight" line**, which is represented by the blue color on the map.

**claude-3-5-sonnet-20240620:** The line that connects Trolskaholm to Solköping is the Bullet line, which is represented by the red line on the map. **This red line runs from Trolskaholm through several stations including Alfheim Central and Singer Park**, before reaching Solköping in the southern part of the map.

**Answer: The Ferry Line**

(a) Map Failure Case

**Explain this image.**

**Answer:** *The inside fur of the cat looks like another cat, so the cat's square root is its own cat-looking fur.*

**gemini-1.5-pro-exp-0827:** The image is a humorous meme playing on the mathematical concept of square roots.

**The Setup:**

* It shows two pictures of the same cat from behind.
* The cat has a distinctive black marking on its back that vaguely resembles a square root symbol (√).

**The Joke:**

* The first image is labeled "Cat".
* The second image is labeled "√Cat".

**The Meaning:**

The joke lies in the visual pun. **The cat's marking looks like a square root symbol,** so the second image humorously implies that the cat has been "square rooted" – as if it were a mathematical equation. It's a silly and absurd concept, which is what makes it funny.

**gpt-4o-2024-05-13:** This image is a humorous visual pun involving a cat.

On the left side, there is a picture of a cat sitting on the ground, viewed from behind. Above this image, the word "Cat" is written.

On the right side, there is a similar picture of the same cat, but the image has been edited to make the cat appear wider and shorter. Above this image, the text "√Cat" (square root of cat) is written.

The joke is that **the cat on the right looks like it has been "squared" (increased in width and decreased in height), playing on the mathematical concept of squaring and square roots.** The cat on the right is meant to humorously represent the "square root" of the cat on the left.

**claude-3-5-sonnet-20240620:** This image is a humorous play on words and mathematics, using pictures of a cat to illustrate the concept.
The top of the image shows two text labels:
1. "Cat" on the left
2. "√Cat" (square root of Cat) on the right
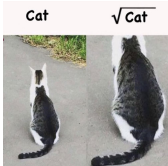Below these labels are two photos of the same cat, presumably:

1. On the left, we see a normal-sized cat sitting on a paved surface. It has a black and white coat.

2. On the right, we see what appears to be the same cat, but it looks elongated or stretched vertically. This creates the illusion that the cat is taller or longer than in the left image.

The joke here is that the **"square root" of Cat (√Cat) is visually represented by a stretched version of the cat, as if taking the square root has "elongated" the animal.** It's a clever visual pun combining mathematical notation with a literal interpretation using cat photos.

This type of humor often resonates well on social media platforms, combining familiar concepts (cats and basic math) in an unexpected and visually amusing way.

(b) Meme Failure Case

**The photo shows a boy's new white sneaker and a girl's black leather shoe. Can you describe what happened?**

**Answer:** *The girl with the black leather shoe stepped on the boy with the white sneaker.*

**gemini-1.5-pro-exp-0827:** Based on the photo, **it appears the boy with the white sneaker has stepped on the girl's black leather shoe.** The white sneaker is positioned over the black shoe, suggesting that it has come down on top of it. This could have happened accidentally in a crowded space, or perhaps during play or a moment of clumsiness. Without more context, it's impossible to know the exact circumstances, but the positioning of the shoes strongly indicates that the white sneaker has stepped on the black shoe.

**gpt-4o-2024-05-13:** The photo shows a **black leather shoe stepping on a white sneaker.** It appears that the girl wearing the black leather shoe has accidentally or intentionally stepped on the boy's new white sneaker. This could have resulted in the white sneaker getting scuffed or dirty.

**claude-3-5-sonnet-20240620:** Based on the image, it appears that the white sneaker and the black leather shoe have come into contact with each other. The white sneaker looks new and clean, while the **black leather shoe seems to be pressing against or stepping on the back part of the white sneaker.** This situation often occurs when someone accidentally steps on the back of another person's shoe while walking closely behind them. It can lead to the heel of the shoe being pushed down or the back of the shoe getting scuffed. The contrast between the pristine white sneaker and the darker leather shoe is quite noticeable in the image. This kind of incident can be frustrating, especially with new shoes, as it might leave marks or damage the shoe's appearance. It's a common mishap that can happen in crowded areas or when people are walking in close proximity to each other.

(c) Shoes Failure Case

Figure 27. **Hard Reasoning Failure Cases (part 1).** The failure cases highlight the model's inability to connect the visual reasoning with language reasoning. While the vision model identifies an object, it is not able to identify the relationship between them correctly.

## (a) Chess Failure Case

**I'm playing as white, what should I do in this position?? I'm confused.**

*Possible Answer: Re7 Rb7 to take the Rook*

**gemini-1.5-pro-exp-0827:** The best move for white in this position is
**Rxe7+**.
**Explanation:**
**1. Rxe7+ forces the black king to capture the rook (Kxe7).**
**2. Qxg7+ follows, checking the black king again.**
3. The black king is forced to move to h8 (Kh8).
4. **Qxh8#** delivers checkmate.
**This is a forced checkmate sequence that wins the game for white.**
**Why other moves are not as good:**
* Moving pawns on the queenside doesn't achieve anything significant and allows black to consolidate their position.
* Developing the knight to f3 doesn't pose an immediate threat and black can easily defend.
* Other rook moves don't lead to a forced checkmate.

**Therefore, Rxe7+ is the only winning move for white in this position.**

**gpt-4o-2024-05-13:** In this position, you have a strong attacking opportunity as White. You can deliver a checkmate in two moves. Here's what you should do:
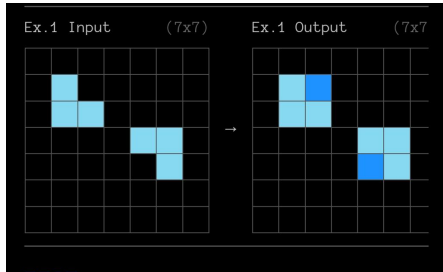
1. **Rd8+** (Rook to d8 check)
- This move forces the Black king to move to e7 because the rook on d8 is **protected by the knight on c5.**

2. **Re8#** (Rook to e8 checkmate)
- The rook on e8 delivers checkmate as the Black king has no legal moves and cannot escape the check.

So, the sequence is:
1. Rd8+ Ke7
2. Re8#

**This results in a checkmate.**

**claude-3-5-sonnet-20240620:** In this chess position, you're playing as White and have a few good options to consider. Here are some suggestions:
1. The most immediate threat is **Black's bishop on f6**, which is attacking **your rook on d4**. You should address this threat first.
2. One strong move would be to capture the **black pawn on e7** with your rook: Rxe7+. This move puts Black's king in check and removes the defender of the f6 bishop.
3. After Black responds to the check (likely by moving the king to f8), you can then capture the bishop on f6 with your bishop: Bxf6.
4. This sequence would result in you winning a pawn and a bishop for the exchange of one of your rooks, which is generally a good trade.
5. Another option, if you want to play more cautiously, is to simply move your rook away from the attack, perhaps to d3 or d2.
6. You could also consider developing your pieces further, such as moving your knight on c5 to a more active square like e6, which would attack Black's rook on b7.

Remember, chess is about long-term strategy as well as immediate tactics. Consider how your move will affect the overall position and try to improve your piece placement while addressing immediate threats.

## (b) ARC Failure Case

**Convert the image into a text based format for both the input and output image, use "[ ]" to represent each square in a grid, "[b]" to represent light blue grids and "[B]" to represent dark blue grids**

**gemini-1.5-pro-exp-0827**  **Output:**
**Input:**
```
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][b][ ][ ][ ][ ][ ]          [ ][b][B][ ][ ][ ][ ]
[ ][b][b][ ][ ][ ][ ]          [ ][b][b][ ][ ][ ][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][ ][ ][b][b][ ][ ]          [ ][ ][ ][B][b][ ][ ]
[ ][ ][ ][ ][b][ ][ ]          [ ][ ][ ][b][b][ ][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
```

**gpt-4o-2024-05-13**

### Ex.1 Input (7×7)          ``` ### Ex.1 Output (7×7)
```
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][b][b][ ][ ][ ][ ]          [ ][b][B][ ][ ][ ][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][ ][ ][ ][b][b][ ]          [ ][ ][ ][ ][b][b][ ]
[ ][ ][ ][ ][ ][ ][b]          [ ][ ][ ][ ][ ][ ][B]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
```

**claude-3-5-sonnet-20240620**

Input (7×7):          Output (7×7):
```
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][b][ ][ ][ ][ ][ ]          [ ][b][B][ ][ ][ ][ ]
[ ][b][b][ ][ ][ ][ ]          [ ][b][b][ ][ ][ ][ ]
[ ][ ][ ][ ][b][ ][ ]          [ ][ ][ ][ ][b][ ][ ]
[ ][ ][ ][b][ ][b][ ]          [ ][ ][ ][b][ ][b][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][B][ ][ ][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
```

**Answer:**
```
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][b][ ][ ][ ][ ][ ]          [ ][b][B][ ][ ][ ][ ]
[ ][b][b][ ][ ][ ][ ]          [ ][b][b][ ][ ][ ][ ]
[ ][ ][ ][b][b][ ]              [ ][ ][ ][b][b][ ]
[ ][ ][ ][ ][b][ ]              [ ][ ][ ][B][b][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
[ ][ ][ ][ ][ ][ ][ ]          [ ][ ][ ][ ][ ][ ][ ]
```

Figure 28. **Hard Reasoning Failure Cases (part 2).** These failure cases highlight the inability for the model to be able to correctly map out a grid-like structure and the various pieces in it.