

CONTEXTUAL: Evaluating Context-Sensitive Text-Rich Visual Reasoning in Large Multimodal Models

Rohan Wadhawan^{*1} Hritik Bansal^{*1} Kai-Wei Chang¹ Nanyun Peng¹

¹University of California Los Angeles

con-textual.github.io

Abstract

Recent advancements in AI have led to the development of large multimodal models (LMMs) capable of processing complex tasks involving joint reasoning over text and visual content in the image (e.g., navigating maps in public places). This paper introduces CONTEXTUAL, a novel benchmark comprising instructions designed explicitly to evaluate LMMs' ability to perform context-sensitive text-rich visual reasoning. CONTEXTUAL emphasizes diverse real-world scenarios (e.g., time-reading, navigation, shopping and more) demanding a deeper understanding of the interactions between textual and visual elements. Our findings reveal a significant performance gap of 30.8% between the best-performing LMM, GPT-4V(ision), and human capabilities using human evaluation indicating substantial room for improvement in context-sensitive text-rich visual reasoning. Notably, while GPT-4V excelled in abstract categories like meme and quote interpretation, its overall performance still lagged behind humans. In addition to human evaluations, we also employed automatic evaluation metrics using GPT-4, uncovering similar trends in performance disparities. We also perform a fine-grained evaluation across diverse visual contexts and provide qualitative analysis which provides a robust framework for future advancements in the LMM design.

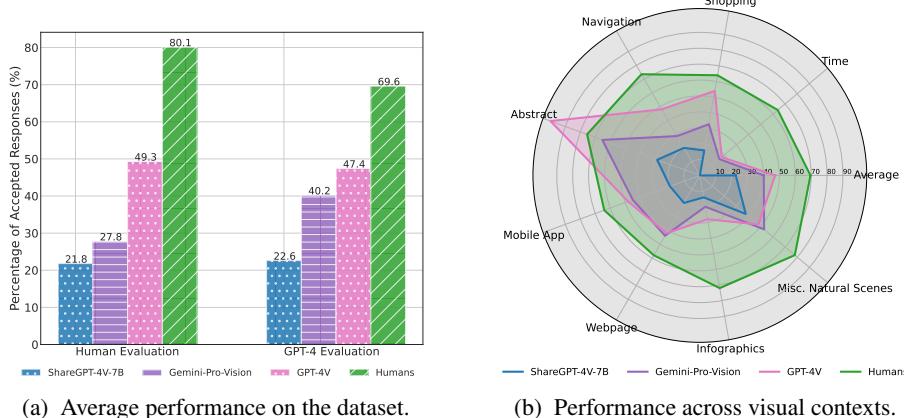


Figure 1: Performance of GPT-4V, Gemini-Pro-Vision, ShareGPT-4V-7B, and Humans on the CONTEXTUAL dataset. (a) Human evaluation and an automatic GPT-4 based evaluation of the response correctness. (b) Fine-grained performance with variation in visual contexts using GPT-4 based evaluation.

^{*}Equal contribution. Contact rwdhawan7@cs.ucla.edu, hbansal@ucla.edu.

1 Introduction

A long-standing goal in AI is to build multimodal assistants that can interact with humans to solve diverse tasks. In this regard, recent development of instruction-tuned large multimodal models (LMMs) have resulted in models that can respond to human instructions, posed as questions or imperative tasks, over images [1, 2, 3, 4, 5]. Naturally, many images in our surroundings have text within them which provides cues for their comprehensive understanding. The ability to reason about the interactions between the text and visual context in the images powers many real-world applications. For example, interpreting text-rich scenes (e.g., navigating maps in public spaces) for assisting the visually impaired, and creative understanding of abstract text-rich images (e.g., memes).

In prior work, numerous datasets have been proposed to access the text-rich visual reasoning of the models such as TextVQA [6], STVQA [6], ESTVQA [7]. These datasets are focused at accessing the OCR capability of the models to *read* the text in the image. As a result, these datasets do not require the model to capture the visual context in the image to answer the question. In Figure 2(a), we highlight a few examples from the ESTVQA dataset. Here, we show that a high accuracy OCR of the images (e.g., ‘W 63 ST...Monk’) has sufficient signal to answer the question (e.g., ‘What is the street number?’). Though accessing the OCR capability is important, these examples do not test the unique potential of the LMMs to jointly reason over the text and visual context in the image.

To address this challenge, we propose **ConTEXTUAL**, a **C**ontext-sensitive **T**ext-rich **V**isual reasoning benchmark consisting of **506** challenging instructions for LMM evaluation. We create a diverse set of instructions on text-rich images with the constraint that they should require context-sensitive joint reasoning over the textual and visual cues in the image. As shown in Figure 2(b), we show a few examples from our dataset. We observe that the instruction (‘Which letters are reflected over the water on the ground?’) cannot be answered even by perfectly capturing the OCR of the text content within the image (e.g., ‘RR Water on Pavement...Planes’). While prior work such as MathVista [8] and MMMU [9] assess vision-language math reasoning and multi-discipline (science, art, history) understanding at college level, our benchmark contains open-ended instructions to test text-rich visual reasoning over daily-life natural or digital scenes (e.g., time reading, shopping, navigation). We summarize our novel contributions in comparison to the related works in Table 1. Each instance



Figure 2: Comparison between the characteristics of the existing datasets (e.g., ESTVQA) and ConTEXTUAL dataset for text-rich visual reasoning. (a) The instances from the prior datasets majorly test the ability of the LMM to read. As a result, their questions can be solved by text-only reasoning over the accurate OCR detection. However, we would want to assess more challenging scenarios with the modern models as they are shown to have enhanced visual perception and reasoning abilities. (b) The instances from ConTEXTUAL are constructed such that they test the ability of the model to capture the context under which text and visual contents in the image interact with each other. Here, text-only reasoning of the detected OCR is not sufficient to perform well.

contains a human-written instruction (question or imperative task), and a human-written ground-truth response. We create ConTEXTUAL dataset to cover **eight** real-world scenarios with text in images:

time reading, shopping, navigation, abstract scenes, mobile application, webpages, infographics, and miscellaneous natural scenes. We present an example from each visual context in Figure 3. The diverse visual nature of these categories enables us to conduct a detailed, nuanced evaluation of the model’s capabilities (§3.4).

We conduct extensive experiments on CONTEXTUAL dataset to assess the reasoning abilities of 13 foundation models (§3.1). This set includes three augmented LLMs setups (e.g., GPT-4 [10] prompted with combinations of image OCR, image layouts, and image captions), two proprietary LLMs (e.g., GPT-4V[2], Gemini-Pro-Vision [3]), and eight open LMMs (e.g., LLaVA-1.5 [11], ShareGPT-4V[12], Idefics [13]). In addition, we perform few-shot experiments for a selected set of models (e.g., Gemini-Pro-Vision, Idefics) to analyze the effect of in-context examples on the model’s performance. Further, we establish a human baseline by asking human annotators to write responses to the dataset instructions. Finally, we perform human and automatic evaluations to assess the correctness of the predicted responses with respect to the ground-truth responses in the dataset (§3.2). Overall, we cover a comprehensive set of models and evaluation metrics to track the progress in the text-rich visual reasoning (§3.3).

Through our human evaluation on randomly selected 280 instances, we find that GPT-4V(ision) is the best performing LMM on the CONTEXTUAL dataset where it achieves 49.3% acceptance rating to its generated responses (Figure 1a). Despite this, the performance lags way behind the human baseline of 80.1% which indicates a large gap in the capabilities of the GPT-4V. In addition, we find that the best performing open-model, ShareGPT-4V-7B, achieves 21.8% rating which indicates that the capabilities of open models are way behind the proprietary models on context-sensitive text-rich visual reasoning (§3.3). Our results highlight that the CONTEXTUAL is a challenging dataset for modern LMMs while humans excel on it.

Since human evaluations are hard to scale and expensive, we also perform automatic evaluation (e.g., GPT-4, GPT-4V, BLEURT [14]) on the complete dataset for all the models (§3.3.1). We present the results with GPT-4 evaluation in Figure 1a along with the fine-grained analysis in Figure 1b. We find that the humans outperform the best performing model GPT-4V by 22.4% highlighting a large gap to fill. Notably, GPT-4V outperforms human performance on the abstract category which involves understanding of visual contexts such as memes and quotes.

Finally, we conduct a qualitative assessment of model predictions, demonstrating a range of performance levels. GPT4V and Gemini-Vision-Pro exhibit superior performance, while open-source LMMs and the Augmented LLM approach lag. Amongst the models, GPT4V shows the highest performance due to its strong visual perception and context-sensitive text-vision reasoning abilities. Further, qualitatively analyzing model responses for instructions from CONTEXTUAL also helps identify issues like hallucination, lack of grounding the instruction to the image, and fine-grained visual perception associated with poor context-sensitive text-vision reasoning (§4).

Table 1: Comparison with related works for evaluating large multimodal models for their text-rich visual reasoning. We abbreviate Context-sensitive as Consens., Generation as Gen.

	Ours	LLaVA[1]	VisIT[15]	TextVQA[6]	STVQA[16]	DUDE[17]	InfographicVQA[18]	SEED [19]
Year Released	2024	2023	2023	2019	2019	2023	2021	2023
Consens. Text-Rich Visual Reasoning	✓	✗	✗	✗	✗	✗	✗	✗
Text in Images	✓	✗	✗	✓	✓	✓	✓	✗
Number of LLM/LMM Models	13	3	10	-	-	9	-	15
Number of Images	506	24	574	28.4K	23K	5K	5.4K	19K
Diverse Image Sources	✓	✗	✓	✗	✓	✓	✗	✗
Question Instructions	✓	✓	✓	✓	✓	✓	✓	✓
Imperative Instructions	✓	✗	✓	✗	✗	✗	✗	✗
Instruction Gen. by Humans	✓	✓	✓	✓	✓	✓	✓	✗
Reference Response Gen. by Humans	✓	✓	✗	✓	✓	✓	✓	✗
Human Evaluation	✓	✗	✓	✓	✗	✓	✓	✗
Automatic Evaluation	✓	✓	✓	✓	✓	✓	✓	✓
Human-Auto Eval. Correlation	✓	✗	✓	✗	✗	✗	✗	✗
Human performance	✓	✗	✗	✓	✗	✓	✓	✗
Absolute Score to Models	✓	✓	✓	✓	✓	✓	✓	✓
Fine-grained Analysis	✓	✗	✓	✗	✗	✓	✗	✓

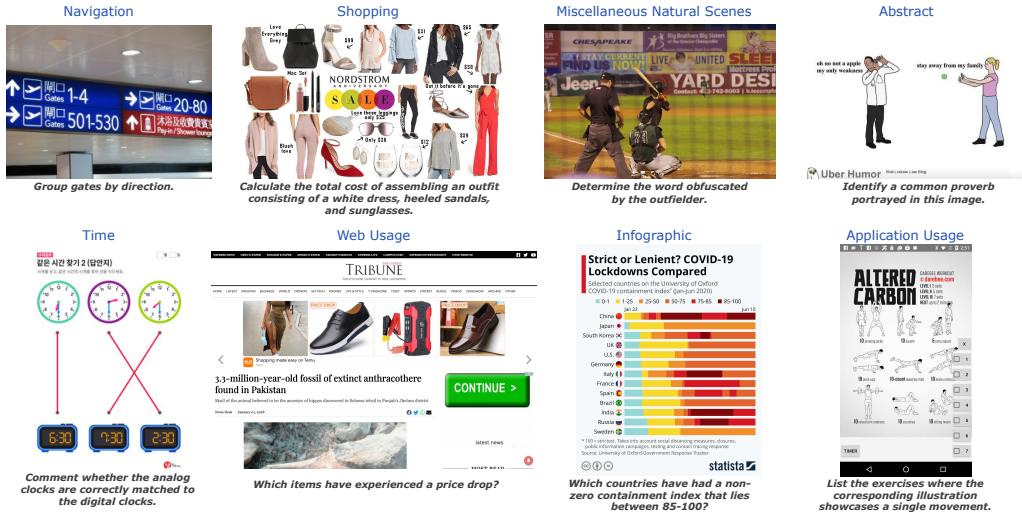


Figure 3: A sample from the 8 visual contexts in CONTEXTUAL that large multimodal models should be able to follow.

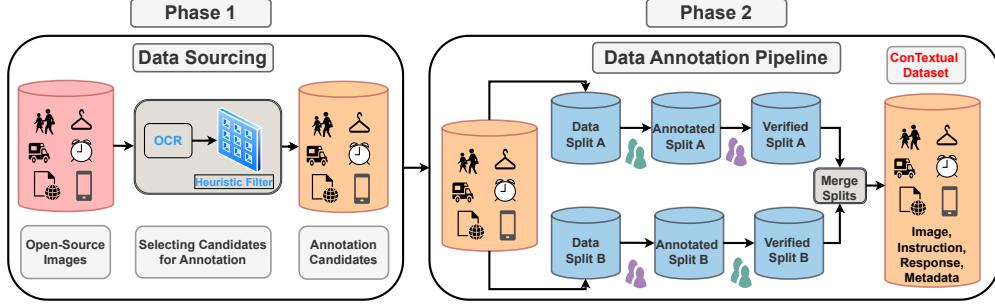


Figure 4: Data Collection Pipeline: (1) ‘Shortlist Images - Utilize manual and automatic methods to shortlist images from the source dataset for annotation. (2) Instruction-Response Creation - Divide the authors into two distinct groups (Group 1 in green and Group 2 represented in purple), assigning each group the responsibility of annotating images corresponding to four categories. (3) Data Sample Validation - Scrutinize the $<image, instruction, response>$ triplets annotated in the prior stage, with one group cross-validating the annotations made by the other group.

2 The CONTEXTUAL Dataset

2.1 Collection Guidelines

We note that there is a notable gap in the existing benchmarks for text-rich images, which primarily evaluate text reading capabilities of LMMs. Our dataset bridges this gap and offers an evaluation framework to test the joint reasoning capabilities of the LMMs over the embedded text and the visual features in the image (Figure 2). Our dataset encompasses a variety of tasks across diverse natural and digital text-rich visual scenarios, thus enabling robust testing.

Broadly, our benchmark follows these key dataset collection guidelines:

- Each sample consists of an $<image, instruction, response>$ triplet, such that the instructions require the models to perform context-sensitive reasoning over the text and visual elements in the image. Specifically, we would avoid creating instructions that could be answered by text-based reasoning (e.g., using LLM) over the detected OCR.

- We aim to cover diverse instructions, including questions and imperative tasks. This ensures that the resulting dataset demonstrates a rich variety of instructions.
- We aim to create instructions of varied complexity. Specifically, they can make extractive instructions that involve extraction of a specific textual or visual elements (e.g., ‘Name the drink with banana flavor.’). In addition, they can make instructions that involve reasoning over the embedded information (e.g., ‘count the number of words in the rendered in the blue color.’).

In this work, we establish a taxonomy by categorizing a dataset into eight distinct visual scenarios, encompassing real-world and digital human interactions. These scenarios include Shopping (e.g., buying groceries, clothing, gadgets), Navigation (various modes of transportation and navigation signage), Time (depicting time and dates using clocks and calendars), Web Usage (websites across different domains), Application Usage (smartphone applications covering education, games, entertainment, etc.), Infographic (visual information on various topics), Abstract (text-rich images illustrating abstract concepts like memes), and Miscellaneous Natural Scenes (miscellaneous human interactions). More details are available in Appendix §A.

Statistic	Number
Total number of samples	506
- Images	506
- Instructions	506
- Responses	506
Root verbs in instructions	79
Visual Scenarios	8
Samples per Scenario	
- Time	50
- Shopping	50
- Navigation	50
- Abstract	50
- Application Usage	50
- Web Usage	50
- Infographic	50
- Miscellaneous Natural Scenes	156
Average Instruction Length	65
Average Response Length	117

Table 2: Key Statistics of CONTEXTUAL

2.2 Data Sources

CONTEXTUAL comprises images sourced from six different sources. Firstly, we obtain images for the *Time*, *Shopping*, *Navigation* & *Abstract* categories from the LAION-5B [20]. Specifically, we use keyword search using CLIP-retrieval UI [21]. A keyword consists of category-specific word(s) + “text” (e.g., clothes text for shopping, airports text for navigation). Some category-specific words we used are: shopping (e.g., grocery, furniture, gadgets, cosmetics, services, clothes), navigation (e.g., street signs, cars, buses, trains, metro, airport, stations, highways, roads), time (e.g., clocks, multiple clocks, digital clocks, timezones, calendars, schedule, watch), and Abstract (e.g., memes, quotes, comic strips, science jokes, math jokes, infographic memes).

Secondly, we source images for the *Application Usage* category from the Rico Dataset [22], which includes 66,000 distinct UI screens originating from 9,300 Android apps across 27 categories, each of these screens featuring visual, textual, structural, and interactive design attributes essential for context-sensitive text-rich visual analysis. Thirdly, we scrape the website links made available by the Open WebText Initiative [23] and collect screenshots for the *Web Usage* category. Open WebText is an open-source effort to reproduce OpenAI’s WebText dataset [24]. Lastly, we acquire images from the test sets of three Visual Question Answering (VQA) datasets, namely InfographicVQA [18],

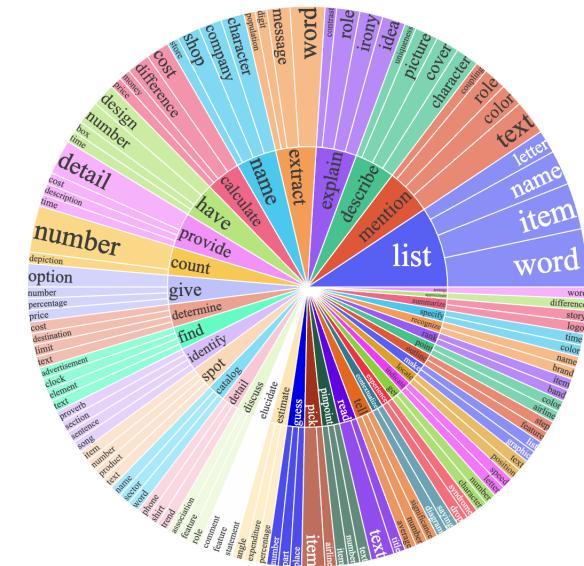


Figure 5: Top 40 Most frequently occurring verbs (inner circle) and their top 4 direct nouns (outer circle) in the instructions.

STVQA [16], and ESTVQA [7], and proceed to annotate them with novel instruction-response pairs. Originally, these datasets consist question-and-answer pairs that primarily focus on text recognition capabilities. However, they offer an opportunity to formulate context-sensitive instructions for the images. Consequently, we reannotate these image instances, categorizing them into two groups: *Infographic*, sourced from the InfographicVQA dataset, and *Miscellaneous Natural Scenes*, sourced from the STVQA and ESTVQA datasets.

2.3 Data Annotation

After image sourcing, we proceed with a three-stage annotation process, shown in Figure 4.

Stage 1: Shortlist Images. In this stage, we shortlist images for instruction-response pair annotation. The images that are categorized under *Time*, *Shopping*, *Navigation*, and *Abstract* undergo manual filtering to guarantee their suitability for annotation. However, for *Application Usage*, *Web Usage*, *Infographic*, and *Miscellaneous Natural Scenes*, we perform heuristic-based filtering. Specifically, we employ a PaddleOCR [25] to detect the text in the image. Subsequently, we select the top 500 images with the highest number of words, a subset of which get annotated in our dataset.

Stage 2: Instruction-Response Creation. Here, we divided the authors into two groups, namely Group 1 and Group 2, each responsible for annotating four specific categories. The authors strictly adhered to the provided annotation guidelines throughout the annotation process.¹

Stage 3: Data Sample Verification. In this final stage, we perform a verification process for each sample annotated in Stage 2. To this end, each group is tasked to review the annotations created by the other group. This ensures adherence to guidelines, and filters out low-quality samples. Finally, we end up with a dataset of **506** instances.

Metadata Annotation Fine-grained metadata facilitates a comprehensive analysis of a benchmark and the performance of models on it. To this end, we provide statistics about the **CONTEXTUAL** benchmark, as shown in Table 2. We visualize each instruction based on its root verb and the direct noun, as shown in Figure 5. We also annotated each sample to determine whether it is an extractive type and whether it requires mathematical reasoning (Appendix §E.1).

2.4 Data Release

CONTEXTUAL comprises 506 samples spanning eight visual categories (refer to Table 2). To facilitate model development, we will release a subset of 100 samples from the 506, as validation set, along with their reference responses, while keeping them hidden for the remaining 406 samples. We ensure that the distribution of validation samples closely mirrors the overall dataset distribution. To achieve this, we randomly select 30 samples from the ‘Miscellaneous Natural Scenes’ category and 10 samples from the remaining categories, maintaining a proportional representation of each category in the validation samples, consistent with the overall benchmark. In this paper, all the results are reported on the entire dataset, unless stated otherwise.

3 Experiments

We conduct a thorough quantitative and qualitative evaluation of the context-sensitive text-rich visual reasoning capabilities of the large multimodal models (LMMs) using CONTEXTUAL. Specifically, we assess 10+ foundation models and establish a human performance baseline on our dataset 3.1. In addition, we perform human evaluation and utilize various automatic evaluation methods to rate the predicted responses from the models 3.2. Further, we provide quantitative results on the model’s performance along with human-automatic correlation analysis in §3.3. Additionally, we perform fine-grained evaluation of the model’s performances under visual contexts, visual scenes and task complexity in §3.4. Finally, we provide qualitative analysis in §4.

¹We observe that MTurk workers found this task time-consuming, leading to annotations that would be hard to accomplish within a limited budget.

3.1 Setup

Augmented LLMs. Prior works [26, 27, 28, 29] have shown that the large language models can solve a wide range of visual reasoning problems by utilizing sensory information from diverse visual tools (e.g., image captions, object detectors models). Since our dataset is focused on text-rich visual reasoning, it is imperative to understand the extent to which a strong LLM GPT-4 can perform on CONTEXTUAL dataset with the OCR information and image captions. To this end, we study this augmented setup under three settings: GPT-4 prompted with (a) vanilla OCR, (b) layout-aware OCR, and (c) combining layout-aware OCR with image captions. We leverage the PP-OCRV4 model of PaddleOCR library [25] for extracting OCR from the images, LATIN prompt [30] inspired OCR text arrangement implementation to maintain layout-awareness in the OCR, and ShareGPT-4V-7B for the dense image captions. We provide the GPT-4 prompt in Appendix §B.

LMMs. We evaluate GPT-4V [2] and Gemini-Pro-Vision [3] that are representative proprietary LMMs that have achieved state-of-the-art on other visual reasoning benchmarks [31]. In addition, we evaluate a wide range of open LMMs including LLaVA-1.5-13B [11], ShareGPT-4V-7B [12], mPLUG-Owl-v2-7B [32, 33], Qwen-VL-Plus [5], InstructBLIP [4], and Idefics-9B [13]. We include LLaVAR [34] and BLIVA [35] as they were introduced for enhanced text-rich visual reasoning.

By default, we assess the performance of the above models under the zero-shot setting. We perform few-shot evaluation (2-shot, 4-shot, 8-shot) for selected models including augmented LLM (GPT-4 with layout-aware OCR and image captions), Gemini-Pro-Vision, and Idefics-9B.

Humans. We also benchmark the performance of humans on our dataset using Amazon Mechanical Turk. The selected annotators that pass an qualification test were asked to write accurate responses for all the instruction-image from the dataset. We provide the screenshot of our annotation interface in Appendix §C.1. We spent \$180 on collecting human predictions on our dataset.

3.2 Evaluation

As discussed in §2, every instance in the CONTEXTUAL dataset consists image, instruction, and a reference response in natural language. The goal of the evaluation is to assess the correctness of the model response for the given instruction for an image.

3.2.1 Human Evaluation

To perform a faithful evaluation of the predicted responses, we ask human annotators sourced from Amazon Mechanical Turk to rate the predicted response quality given the image, instruction, and reference response from our dataset. First, we sample 280 instances from the dataset randomly from the CONTEXTUAL dataset. Second, we collect the model responses for these instances from augmented LLM (GPT-4 with layout-aware OCR and image captions), GPT-4V, Gemini-Pro-Vision, LLaVA-1.5-13B, ShareGPT-4V-7B, and humans. In total, we have 1680 predicted responses from models and humans. Third, we show each model response, without revealing the model identity, to three human annotators independently. Specifically, the human annotators are asked to decide the predicted response is acceptable given the reference response, instruction and image from the dataset. Finally, we report the acceptance rating (0-100 in percentage) of the responses using the majority vote among the three annotator as the final decision for a model prediction. We provide the screenshot of our annotation interface in Appendix C.2. We spent \$1000 in acquiring human judgments.

3.2.2 Automatic Evaluation

While human evaluation acts as a gold standard, it is hard to scale since it is expensive and time-taking. Since our dataset uniquely provides reference response for each instruction, we utilize test a wide range of reference-guided automatic evaluation methods. Specifically, these include (a) prompting an LLM GPT-4 with the instruction, reference response and predicted response, (b) prompting an LMM GPT-4V with the image, instruction, reference response and predicted response, (c) and other text generation methods like BLEURT [14], Rouge-L [36] and BERTScore [37] that assess the similarity between the reference response and predicted response. Specifically, GPT-4 and GPT-4V are prompted to provide their judgement on the predicted response, same as human evaluation. We

present the prompt for GPT-4 based evaluation in Appendix §D. However, the other text generation methods provide a continuous score 0-1 which is scaled to 0-100.

Through our automatic evaluation methods, we evaluate all the model responses on the entire dataset. Subsequently, we conduct a correlation analysis between human and automated methods, utilizing the same 1,680 responses from the human evaluation, to assess the efficacy of the automated approaches (§3.3.1). Finally, we utilize the GPT-4 automatic evaluation, that achieves the highest correlation with human judgments, for large-scale evaluation of all the models on the complete dataset (§3.4).

Table 3: Comparison in the performance of various foundation models (augmented LLM and LMMs) and humans on the CONTEXTUAL dataset. We report the response acceptance rating using human evaluation, automatic GPT-4 and GPT-4V based evaluation. In addition, we report standard text generation quality assessment metrics including BLEURT, Rouge-L, and BERTScore. We find that humans outperform existing models on our dataset. The best performing LMM is GPT-4V is behind human performance by 30%. The best performance in a column is highlighted in **BLACK** while the second best performance is highlighted in UNDERLINE.

	Humans	GPT-4	GPT-4V	BLEURT	Rouge-L	BERTScore
GPT-4 w/ Layout-aware OCR + Caption	17.2	22.2	17.6	41.3	22.5	53.9
GPT-4V [2]	<u>49.3</u>	<u>47.4</u>	<u>45.0</u>	<u>45.3</u>	17.3	52.5
Gemini-Pro-Vision [3]	27.8	40.2	37.1	42.5	<u>30.1</u>	<u>58.4</u>
LLaVA-1.5-13B [11]	17.2	20.6	17.5	43.6	21.7	54.8
ShareGPT-4V-7B [12]	21.8	22.6	20.6	44.5	23.3	55.8
Humans	80.1	69.6	68.6	47.4	33.6	59.8

3.3 Results

We compare the performance of augmented LLM, LMMs, and humans on CONTEXTUAL using human and automatic evaluation in Table 3. Through our human evaluations, we find that the humans perform the best on the dataset with the response acceptance rating of 80.1%. In addition, we observe that the GPT-4V achieves the highest acceptance rating of 49.3% in comparison with all the other models. However, this rating is quite far from the human performance which indicates that our task is quite challenging for the state-of-the-art LMMs while humans are good at it. We find that the GPT-4V outperforms Gemini-Pro-Vision by 22% highlighting a large gap in the models text-rich visual reasoning capabilities. Further, we find that augmented LLM approach achieves a very low rating of 17.2% which indicates that the dataset instances cannot be solved without precise visual perception. Interestingly, we observe that the open-models such as LLaVA-1.5-13B and ShareGPT-4V-7B achieve poor acceptance ratings through human evaluations which indicates the presence of a large gap in their capabilities from proprietary models. This might be attributed to the differences in the model capacity, along with the scale and quality of the pretraining data.

As human evaluation is not scalable, we perform automatic evaluation of the model responses on the entire dataset. In Table 3, we find that the ratings of the human responses outperforms the ones from GPT-4V by 22.2% and 23.6% using GPT-4 and GPT-4V evaluation. Alike human evaluation, automatic evaluation with GPT-4 and GPT-4V highlights that the human performance on the CONTEXTUAL dataset is way higher than the best performing LMM. Interestingly, the gap between the performance GPT-4V and Gemini-Pro-Vision is 7.2%. In addition, we still observe a large gap in performance of the proprietary models and open LMMs. We perform fine-grained evaluation to understand the gaps in model capabilities along the various quality dimensions in §3.4.

In addition, we find that the BLEURT scores for the humans is the highest while GPT-4V achieves the highest score among the LMMs. Interestingly, we find that the open models (LLaVA-1.5, ShareGPT-4V) achieve a higher BLEURT score than Gemini-Pro-Vision. We observe similar counter-intuitive trends in our Rouge-L and BERTScore based automatic evaluations. For instance, Rouge-L and BERTScore rank open models better than GPT-4V despite considering the human responses to be the best. This counter-intuitive observation might be attributed to the sensitivity of these methods to the differences in lexical variations in the reference and predicted responses [14].

3.3.1 Correlation Analysis of Automatic Evaluation Methods

We measure the correlation between the candidate automatic metrics and human judgments using ROC-AUC and spearman correlation in Table 4. Specifically, the human judgments are considered as gold standard where we assign ‘0’ to unaccepted responses to the instructions and ‘1’ to the

Table 4: Comparison of the human and automatic evaluation metric using ROC-AUC and spearman correlation. We find that the GPT-4 and GPT-4V based evaluation correlate the most with the humans in both the methods.

	GPT-4	GPT-4V	BLEURT	RougeL	BERTScore
ROC-AUC	85.9	83.9	72.9	67.6	66.8
Spearman Correlation	0.71	0.68	0.38	0.29	0.28

accepted responses. We find that GPT-4 based evaluation achieves the highest ROC-AUC of 85.9 and spearman correlation of 0.71 amongst all the automatic evaluation metrics. In addition, we observe that GPT-4V also achieves a high correlation with the human judgments which is close to GPT-4. The differences in the judgments of GPT-4 and GPT-4V can be attributed to the difference in their respective inputs. Specifically, GPT-4 bases its judgments on the given instruction and the reference response, whereas GPT-4V, with access to an input image, may potentially be biased. This access might lead GPT-4V to overlook the reference response and depend on the visual cues from the input image for making judgments in some cases. Finally, we observe that standard text generation metrics achieve a poor ROC-AUC and Spearman correlation in comparison to GPT-4 metrics. This corroborates the findings from the prior research [15] that shows GPT-4 evaluation outperforms standard text generation metrics. As a result, we utilize GPT-4 for automatically evaluate the quality of the predicted responses on the CONTEXTUAL dataset at scale.

3.3.2 Few-Shot Setting

Here, we compare the performance of the foundation models on CONTEXTUAL using GPT-4 evaluation with under the few-shot settings in Figure 6. Specifically, we perform zero-shot, two-shot, four-shot, and eight-shot evaluation for augmented LLM (GPT-4 prompted w/ layout aware OCR and image caption), Gemini-Pro-Vision, and Idefics-9B. We select in-context examples at random from our dataset and evaluate the models on the remaining instances.

In our experiments, we find that the performance of all the models does not change drastically with in-context examples. Specifically, we observe that Gemini-Pro-Vision response acceptance rating decreases by 5% in the two-shot setting as compared to the zero-shot setting, and, increases monotonically from two-shot to eight-shots. In addition, we observe that the performance improvements stagnate for Idefics-9B after the four in-context examples. Recent studies highlight the instability and sensitivity of LMMs in few-shot settings [38]. For instance, a significant accuracy drop was observed in models like InstructBLIP in four-shot setting, especially in tasks requiring commonsense reasoning. Overall, we highlight that providing few-shot examples does not elicit context-sensitive text-rich visual reasoning in the foundation models.

3.4 Fine-Grained Evaluation

We compare the fine-grained performance of a wide range of foundation models across different visual contexts using GPT-4 evaluation in Table 5. In our experiments, we find that GPT-4V outshines the baseline models in almost all categories. We observe that the sole exceptions are web usage and miscellaneous natural scenes contexts, where Gemini-Pro-Vision holds the lead. Notably, GPT-4V outperforms humans on reasoning over the abstract category, highlighting that it may have been tuned to reason over a lot of memes and quotes data. In addition, we observe that all the models struggle the most in the time category while humans ace it, a skill which is could be hard to learn from the training data. After time reading, the proprietary LMMs underperform on the infographics category which consists reasoning over data visualizations. Prior work [8, 39] has shown that the existing LMMs underperform humans in reasoning over charts.

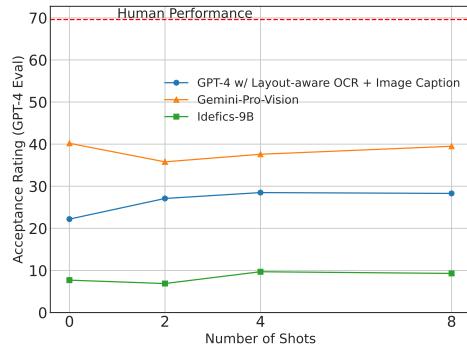


Figure 6: Few-shot performance on the CONTEXUAL dataset.

Table 5: Fine-grained comparison in the performance of the foundation models (augmented LLM and LMMs) and humans on the CONTEXTUAL dataset using GPT-4 evaluation. We abbreviate the average response acceptance rating as Avg., Shopping as Shop., Navigation as Nav., Abstract as Abs., Application usage as App., Web Usage as Web, Infographics as Info., Miscellaneous natural scenes as NS. We find that the GPT-4V outperforms all the model baselines on most of the categories while Gemini-Pro-Vision is the best on Web usage and natural scenes. The best performance in a column is highlighted in **BLACK** while the second best performance is highlighted by UNDERLINE.

MODELS	Avg.	Time	Shop.	Nav.	Abs.	App.	Web.	Info.	Misc. NS.
<i>Augmented Large Language Models</i>									
GPT-4 w/ OCR	15.9	4.0	10.0	14.0	30.6	8.0	16.0	28.6	16.9
GPT-4 w/ Layout-aware OCR	18.2	8.0	20.0	18.0	34.7	10.0	16.0	16.0	20.7
GPT-4 w/ Layout-aware OCR + Caption	22.2	6.0	16.0	24.0	57.1	14.0	18.0	8.0	27.3
<i>Large Multimodal Models</i>									
GPT-4V [2]	47.4	18.0	54.0	48.0	100.0	48.0	42.0	28.0	48.0
Gemini-Pro-Vision [3]	40.2	16.0	32.7	28.6	65.3	44.9	<u>43.8</u>	20.0	<u>52.8</u>
LLaVA-1.5B-13B [11]	20.8	4.0	10.0	18.0	44.9	16.0	26.0	4.0	29.7
ShareGPT-4V-7B [12]	22.6	0.0	16.0	20.0	28.6	20.0	20.0	14.0	37.7
mPLUG-Owl-v2-7B [32]	18.6	4.0	8.0	24.0	32.7	20.0	10.0	12.0	26.0
Qwen-VL-Plus [5]	21.8	4.0	20.0	24.0	53.1	6.0	18.0	14.0	27.3
InstructBLIP-Vicuna-7B [4]	9.7	2.0	4.0	16.0	20.0	6.0	12.0	2.1	12.0
Idefics-9B [13]	7.7	4.0	2.0	12.0	12.0	0.0	6.0	2.0	13.3
LLaVAR [34]	14.9	10.0	16.0	6.0	44.9	8.0	10.0	6.0	16.7
BLIVA [35]	10.3	2.0	4.0	14.0	24.5	4.0	8.0	4.0	14.7
Humans	69.6	64.0	64.0	73.5	<u>75.5</u>	64.0	58.0	72.0	78.0

Further, we observe that the best performing open models like LLaVA-1.5-13B and ShareGPT-4V-7B achieve an acceptance rating above 25% only in the context of abstract and natural scenes while they struggle with time, shopping, navigation, app (mobile) usage, web usage, and infographics. This can be attributed to the lack of diverse visual contexts in their training data. For instance, COCO dataset [40] used for vision-language alignment in the open models, predominantly comprises natural scenes. In addition, we observe that the open models that were specifically introduced for text-rich visual reasoning like LLaVAR and BLIVA falter on CONTEXTUAL dataset. This indicates that these models cannot reason when the instruction requires them jointly over the text content and visual context in the image. We perform additional fine-grained evaluation in Appendix §E. Overall, our fine-grained analysis aids in identifying the gaps in the existing models which would inspire the development of next generation LMMs.

4 Qualitative Examples

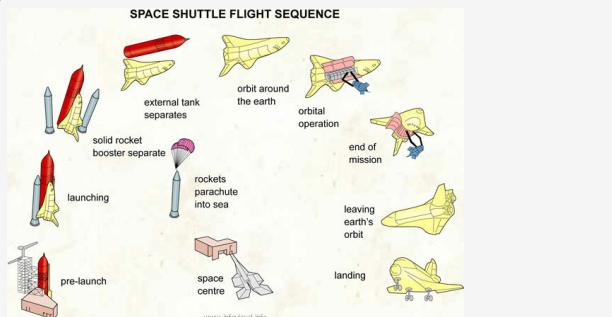
GPT-4V vs Humans. GPT-4V performs well on CONTEXTUAL but still has a gap of **30.8%** according to human evaluation and **22.2%** according to GPT-4 eval. To better understand this gap, we qualitatively analyze GPT-4V’s results compared to human performance. For example, in Figure 7, we see an instance where GPT-4V provides an incorrect answer. In this case, the model is asked to identify the colors of different parts of a space launch vehicle - space shuttle, external tank, and rocket thrusters. GPT-4V makes errors in color predictions but can accurately infer the diagram’s information, revealing a lack of precise visual perception. It wrongly assumes the space shuttle is white and hallucinates the red tips of the rockets. We provide more examples in Appendix §F (Figures 25, 29, 33, 34, 48, 51, 52), highlights that GPT-4V’s core issue lies in fine-grained perception coupled with a bias for prior visual knowledge (e.g., assuming the space shuttle is white even when it’s represented as yellow). Similar analysis was presented in the prior work [41] where GPT-4V fails on the perturbed versions of common visual illusions.

GPT-4V vs. Open LMMs and Augmented LLM. We also qualitatively compare the performance of GPT-4V against open-source LMMs and the Augmented LLM approach. To this end, we analyze the best performing open-source LMM and Augmented LLM approach, that is, ShareGPT-4V-7B and GPT-4 w/ Layout-aware OCR + Caption, respectively, using an example illustrated in Figure 8. In this example, the models are instructed to provide the cost of hiring the illustrated washing machine for nine months. To accurately respond to the instruction, a model must first understand which machine is depicted in the image and then, based on it, extract the appropriate cost for the period mentioned in the instruction. We observe that GPT-4V can correctly identify the machine and

also identify the correct amount. However, ShareGPT-4V-7B and the Augmented LLM approach fail to do so.

The open-source LMM correctly determines the time period but fails to identify the correct machine, highlighting its weakness in context-sensitive reasoning and fine-grained visual perception. This qualitative example is consistent with other instances where open-source LMMs demonstrated inferior context-sensitive reasoning to GPT-4V (Appendix §4 for more examples). They also show a higher tendency to hallucinate, possibly due to a lack of grounding the instruction to the given image or reliance on prior knowledge. We refer to Figure 26,27,30,31,38,47 for more examples. In this example, the augmented LLM approach mispredicts the machine as top load and produces a result similar to that of the LMM. This can be attributed to its inability to reason over text and image jointly (refer to Appendix §4 for more examples).

In summary, GPT-4V outperformed open-source LMMs and Augmented LLMs on CONTEXTUAL, showcasing superior context-sensitive text-vision reasoning. We qualitative analysis helps in identifying issues like hallucination, lack of grounding the instruction to the image, and fine-grained visual perception associated with poor context-sensitive text-vision reasoning.



Instruction: Specify the colors of the rocket, space shuttle, and external tank.

Reference Response: The space shuttle is depicted in yellow, the rockets are depicted in grey, and the external tank is depicted in red.

Human Response: The rocket is gray. The space shuttle is yellow. The external tank is red.

GPT-4V Response: In the image, the solid rocket boosters of the **space shuttle are colored in red and light gray, with the red at the top and the gray on the bottom. The external tank is colored orange. The main body of the space shuttle, also known as the orbiter, is primarily white with black on its nose and wing edges. There is also a bit of red on the tail section.**

Summarized Reasoning: During the pre-launch, all components are vertical. The solid rocket boosters detach when the shuttle ascends. Once the shuttle reaches space's edge, the external tank is released. Throughout these phases, the shuttle's colors remain consistent, i.e., white for the main body that is a typical color for it. External tank is orange, which is its natural insulation color. The solid rocket boosters are white with red tops.



Washing Machine on Hire		
Months	Top Load	Front Load
3 months	₹ 2400	₹ 3600
6 months	₹ 3600	₹ 6000
9 months	₹ 4500	₹ 7200
12 months	₹ 6000	₹ 9600

Instruction: Provide the cost of hiring the machine type depicted in the image for nine months.

Reference Response: The machine shown in the image is a front-loading machine, and the rental fee for a period of 9 months amounts to INR 7200.

GPT-4V Response: The image displays a front load washing machine along with a price list for hiring different types of washing machines. For hiring a front load washing machine for a period of 9 months, the cost would be **INR 7,200** according to the given price list.

ShareGPT-4V-7B Response: The cost of hiring the washing machine for nine months is **4500 rupees**.

GPT-4 w/ Layout-aware OCR + Caption Response: The cost of hiring the depicted machine type (Top Load) for nine months is **4500 units (currency not specified in the provided text)**

Figure 7: In this instance, GPT-4V provides an incorrect response to the instruction, despite its logical reasoning. The use of green indicates responses that match the reference, while red highlights errors in the responses. Additionally, a *Summarized Reasoning* is provided to outline the rationale used by GPT-4V to arrive at its answer.

Figure 8: In this example, GPT-4V correctly responds to the instruction. However, ShareGPT-4V-7B (best performing open-source LMM) and GPT-4 w/ Layout-aware OCR + Caption (Augmented LLM) produce a wrong response, due to lack of joint reasoning over text and image.

5 Related Work

Text-Rich Image Understanding. Recently, there has been a grown interest in understanding the interactions between the text and visual elements in the image [42, 43]. To track the progress of the models in this field, several datasets were introduced like OCRVQA [44], TextVQA [6], DocVQA [45], STVQA [16], ESTVQA [7]. These datasets majorly focus on the ability of the models to accurately read the text in the documents or natural scene images. In addition, there are datasets such as InfographicsVQA [18] and ChartQA [46] that assess model’s ability to reason jointly about the text and visual elements in the image. However, they are focused on specific domains (e.g., infographics) and limited to question answering. Prior work [47] evaluates the ability of LMMs to perform accurate OCR on the existing datasets. In comparison, we propose a new CONTEXTUAL dataset, comprising a wide range of visual contexts, instruction types (questions and imperative tasks), that aims to test the LMM’s ability to perform precise visual perception and complex reasoning over the visual and text elements of the image.

Vision Language Reasoning Benchmarks. Having high-quality datasets is essential to assess the progress of the fields towards building high utility models for the real-world. Traditionally, vision-language learning has focused on tasks such as visual question answering [48, 31] or image captioning [49, 40] where the model primarily needs to understand the key objects and their relations. In addition, there are several datasets and benchmarks that evaluate specific skills of the LMMs including math skills [8], world knowledge [9], and grade school science diagrams [50]. In comparison, our CONTEXTUAL dataset introduces and tests the ability of the large multimodal models to perform context-sensitive text-rich visual reasoning for real-world scenarios. These include time reading, navigation and transportation in public spaces, meme and quote understanding, and shopping etc.

Large Multimodal Models. Recently, there has been a huge surge of vision-language generative models with instruction-following capabilities, termed as large multimodal models (LMMs). These include proprietary models such as GPT-4V [2] and Gemini-Pro-Vision [3]. These models have achieved state-of-the-art performance on the traditional vision-language models. In the open space, the models include LLaVA [1], mPLUG-Owl [32], OpenFlamingo [51], Idefics [13], LLaMA-Adapter [52], Idefics [13]. In addition, there are a class of LMMs that focus on enhanced text-rich visual reasoning capabilities including LLaVAR [34] and BLIVA [35]. In this work, we compare the performance of LMMs on the CONTEXTUAL dataset. We find that the text-rich visual reasoning capabilities of the proprietary models is way superior than the open models. We also include fine-grained analysis to understand the gaps in the model performance across different visual contexts.

6 Conclusion

In this work, we introduce CONTEXTUAL, a dataset for evaluating the text-rich visual reasoning in large multimodal models. Going beyond the prior efforts that focus primarily on the testing the reading skills in the visual contexts, we create novel and challenging instructions from scratch that would require the models to capture the context in which the text is presented in an image. In addition, we perform ask humans to solve our dataset and also use human annotators for model response evaluation. In our experiments, we find that the modern LMMs (proprietary and open models) struggle to perform on our dataset while humans are good at it. We also perform fine-grained evaluation and qualitative analysis to pinpoint the gaps in the model capabilities. In summary, our dataset paves a path for assessing the progress on reasoning over text-rich images, a domain with significant real-world applications.

7 Acknowledgement

This material is based on research supported by the ECOLE program under Cooperative Agreement HR00112390060 with the US Defense Advanced Research Projects Agency (DARPA). We also want to thank the UCLA-Amazon Science Hub for Humanity and Artificial Intelligence. Hritik Bansal is supported in part by AFOSR MURI grant FA9550-22-1-0380.

References

- [1] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.
- [2] OpenAI. Gpt-4v(ision) system card, 2023b. <https://openai.com/research/gpt-4v-system-card>, 2023.
- [3] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [4] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- [5] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- [6] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019.
- [7] Xinyu Wang, Yuliang Liu, Chunhua Shen, Chun Chet Ng, Canjie Luo, Lianwen Jin, Chee Seng Chan, Anton van den Hengel, and Liangwei Wang. On the general value of evidence, and bilingual scene-text visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10126–10135, 2020.
- [8] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- [9] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.
- [10] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023a, 2023.
- [11] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.
- [12] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- [13] HuggingFace. Introducing idefics: An open reproduction of state-of-the-art visual language model. <https://huggingface.co/blog/idefics>, 2023.
- [14] Thibault Sellam, Dipanjan Das, and Ankur P Parikh. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*, 2020.
- [15] Yonatan Bitton, Hritik Bansal, Jack Hessel, Rulin Shao, Wanrong Zhu, Anas Awadalla, Josh Gardner, Rohan Taori, and Ludwig Schimdt. Visit-bench: A benchmark for vision-language instruction following inspired by real-world use. *arXiv preprint arXiv:2308.06595*, 2023.
- [16] Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluis Gomez, Marçal Rusinol, Ernest Valveny, CV Jawahar, and Dimosthenis Karatzas. Scene text visual question answering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4291–4301, 2019.
- [17] Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann, Michał Pietruszka, Paweł Joziak, Rafal Powalski, Dawid Jurkiewicz, Mickaël Coustaty, Bertrand Anckaert, Ernest Valveny, et al. Document understanding dataset and evaluation (dude). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19528–19540, 2023.

- [18] Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1697–1706, 2022.
- [19] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023.
- [20] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.
- [21] Romain Beaumont. Clip retrieval: Easily compute clip embeddings and build a clip retrieval system with them. <https://github.com/rom1504/clip-retrieval>, 2022.
- [22] Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hirschman, Daniel Afergan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th annual ACM symposium on user interface software and technology*, pages 845–854, 2017.
- [23] Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.
- [24] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [25] paddlepaddle. Paddleocr: Multilingual ocr toolkit based on paddlepaddle. <https://github.com/PaddlePaddle/PaddleOCR>, 2023.
- [26] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. *arXiv preprint arXiv:2304.09842*, 2023.
- [27] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- [28] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. *arXiv preprint arXiv:2303.08128*, 2023.
- [29] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14953–14962, 2023.
- [30] Wenjin Wang, Yunhao Li, Yixin Ou, and Yin Zhang. Layout and task aware instruction prompt for zero-shot document image question answering. *arXiv preprint arXiv:2306.00526*, 2023.
- [31] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [32] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- [33] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.04257*, 2023.
- [34] Yanzhe Zhang, Ruiyi Zhang, Juxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. Llavar: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023.

- [35] Wenbo Hu, Yifan Xu, Y Li, W Li, Z Chen, and Z Tu. Bliva: A simple multimodal llm for better handling of text-rich visual questions. *arXiv preprint arXiv:2308.09936*, 2023.
- [36] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [37] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- [38] Yunxin Li, Longyue Wang, Baotian Hu, Xinyu Chen, Wanqi Zhong, Chenyang Lyu, and Min Zhang. A comprehensive evaluation of gpt-4v on knowledge-intensive visual question answering. *arXiv preprint arXiv:2311.07536*, 2023.
- [39] Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq Joty. Unichart: A universal vision-language pretrained model for chart comprehension and reasoning. *arXiv preprint arXiv:2305.14761*, 2023.
- [40] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- [41] Tianrui Guan, Fuxiao Liu, Xiyang Wu Ruiqi Xian Zongxia Li, Xiaoyu Liu Xijun Wang, Lichang Chen Furong Huang Yaser Yacoob, and Dinesh Manocha Tianyi Zhou. Hallusionbench: An advanced diagnostic suite for entangled language hallucination & visual illusion in large vision-language models. *arXiv e-prints*, pages arXiv–2310, 2023.
- [42] Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *International Conference on Machine Learning*, pages 18893–18912. PMLR, 2023.
- [43] Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. Layoutlm: Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1192–1200, 2020.
- [44] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *2019 international conference on document analysis and recognition (ICDAR)*, pages 947–952. IEEE, 2019.
- [45] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021.
- [46] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- [47] Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, et al. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023.
- [48] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.
- [49] 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*, pages 3608–3617, 2018.

- [50] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV* 14, pages 235–251. Springer, 2016.
- [51] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- [52] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. Llama-adapter v2: Parameter-efficient visual instruction model, 2023.

A Dataset Details

A.1 Visual Scenarios Description

In this section, we outline the constituent elements that make up each visual scenario, as illustrated in Table 6.

Category	Description
Shopping	Purchasing groceries, clothes, furniture, gadgets, cosmetics, services, and miscellaneous products.
Navigation	Different modes of transportation - passenger vehicles, trucks, buses, trains, and airplanes, and navigation signage - streets, roadways, bus stations, train stations, and airports.
Time	Items showcasing time and dates, including analog clocks, digital clocks, multi-clock setups, calendars, and other miscellaneous time-viewing setups.
Web Usage	Websites across a variety of domains, like news articles, blogs, sports, and e-commerce
App Usage	Smartphone applications on education, productivity, games, lifestyle, entertainment, news, etc.
Infographic	Infographics on local and global information spanning domains of health, sports, education, natural resources, technology, etc.
Abstract	Memes, comic strips, and other abstract concepts illustrated through text-rich images.
Miscellaneous	Miscellaneous human interactions do not fall into the previous categories.
Natural Scenes	

Table 6: Descriptions of the eight visual scenarios in **CONTEXTUAL**.

A.2 Visual Scenarios Examples

In this section, we provide examples of each visual category in **CONTEXTUAL**.

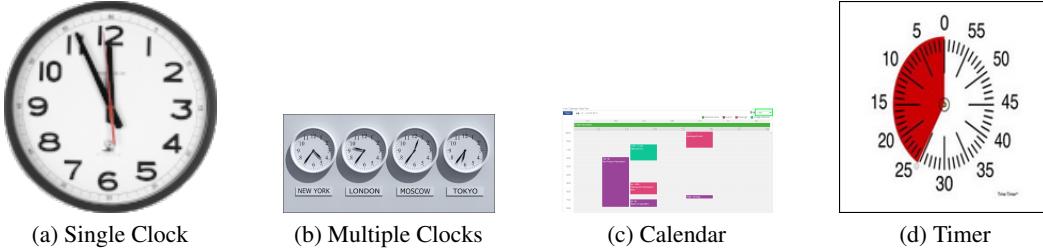


Figure 9: Examples of the *Time* visual Scenario



Figure 10: Examples of the *Shopping* visual scenario

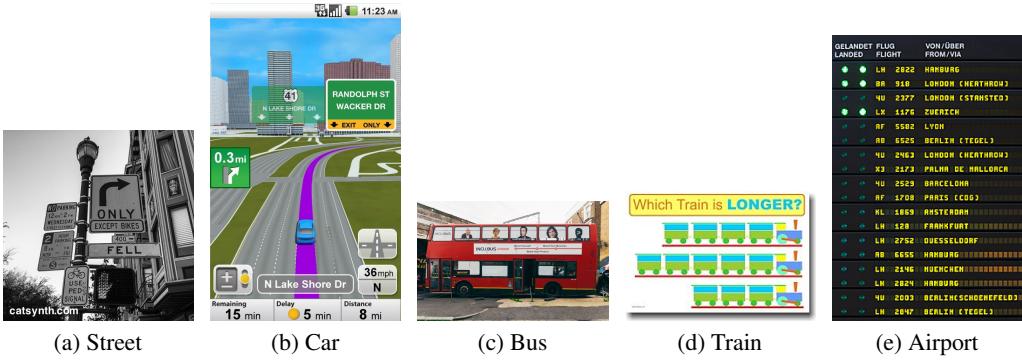


Figure 11: Examples of the *Navigation* visual scenario

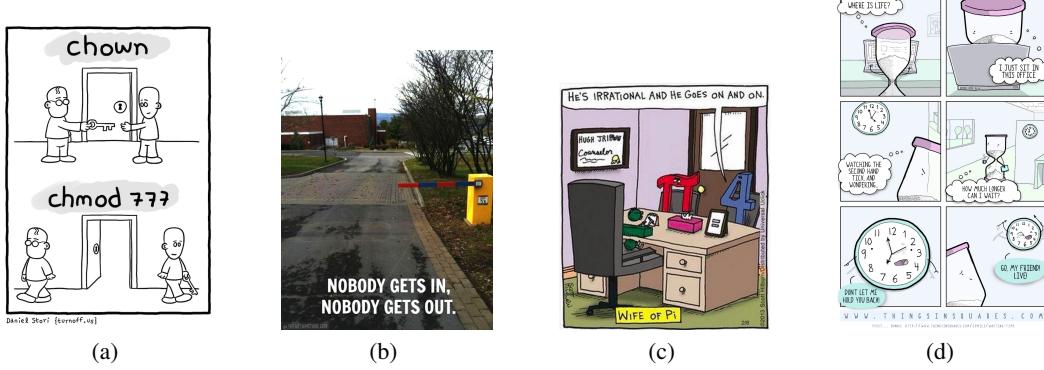


Figure 12: Examples of the *Abstract* visual scenario

B Augmented LLM Prompt

In this section, we discuss the design and elaborate on the prompts employed for the Augmented LLM approach (illustrated in Figure 17, 18, 19). We describe the three distinct prompt formats utilized, each differing in the extent of visual information presented. These formats encompass simple OCR of the image, OCR of the image arranged in the layout it appears in the image, and OCR presented in a layout format along with a comprehensive image caption. We prompt GPT4 with the above templates that does not take the image as input. However, the image is included in the illustration for reference purposes.

C Human Annotation Screenshots

C.1 Human Performance Screenshot

We present the screenshot of the user interface used for acquiring human responses on the CONTEX-TUAL dataset in Figure 20.

C.2 Human Evaluation Screenshot

We present the screenshot of the user interface used for human evaluation in Figure 21.

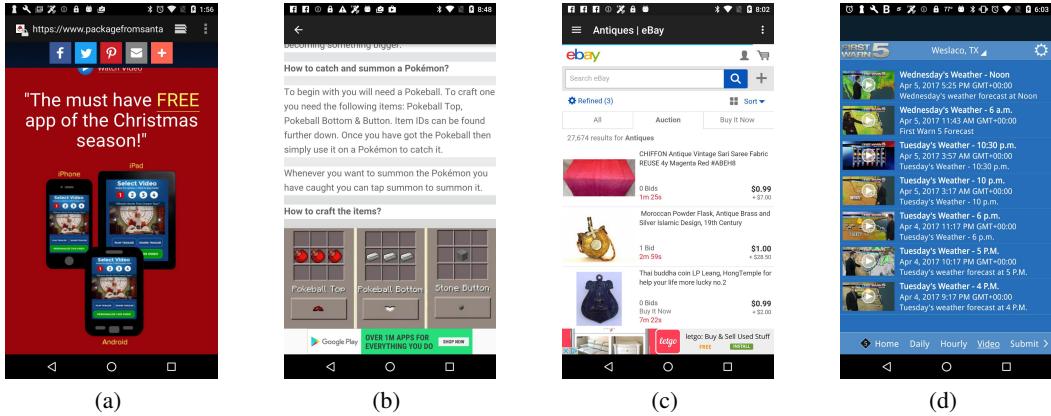


Figure 13: Examples of the *Mobile Usage* visual scenario

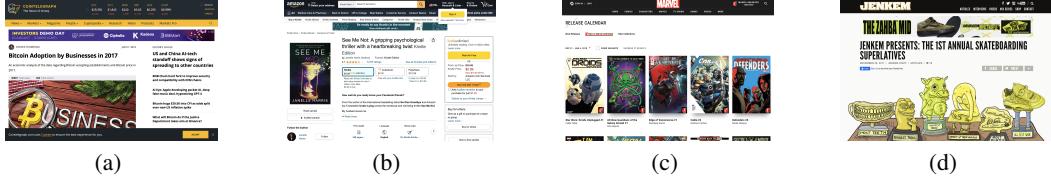


Figure 14: Examples of the *Web Usage* visual scenario

D GPT-4 Evaluation Prompt

E Additional Fine-grained Evaluation

E.1 Types of Tasks

We compare the performance of the foundation models with varying types of tasks in Figure 22a using GPT-4 evaluation. Specifically, we assess the quality of the responses when the instructions require the models to extract text or visual elements in the image (e.g., *List the exercises where the corresponding illustration showcases a single movement.*). There are 285 such instances in the CONTEXTUAL dataset. While these tasks require complex perception and reasoning abilities, they do not require additional operations on top of the information already presented in the image explicitly. We observe that the humans achieve 75.1% on such instructions while the proprietary models GPT-4V and Gemini-Pro-Vision achieve 38.7% and 40.1%, respectively. This indicates that humans are very good at identifying the key information that needs to be extracted to respond to the instructions.

In addition, we assess the responses when the instructions require the models to go beyond information extraction, and perform math reasoning for the instruction (e.g., *What is the total price for the two cars listed here?*). There are 66 instances in the CONTEXTUAL dataset. We find that humans achieve 62.1% on such tasks while the proprietary models GPT-4V achieve 27.3%, again highlighting at the large gap in their math reasoning.

E.2 Visual Scenes

We compare the performance of the foundation models with varying visual scenes (e.g., natural scenes and digital scenes) in Figure 22b. Majorly, shopping, navigation, and misc. natural scenes constitute natural scenes, and web usage, mobile usage, abstract, infographics and time reading constitute digital scenes. We find that humans achieve the highest performance in both the visual scenes i.e., 72.5% and 66.2% on natural scenes and digital scenes, respectively. In addition, we observe that GPT-4V achieves 47.8% and 48.2% on natural and digital scenes, respectively. Interestingly, we find that Gemini-Pro-Vision and ShareGPT-4V-7B achieve higher performance on the natural scenes than

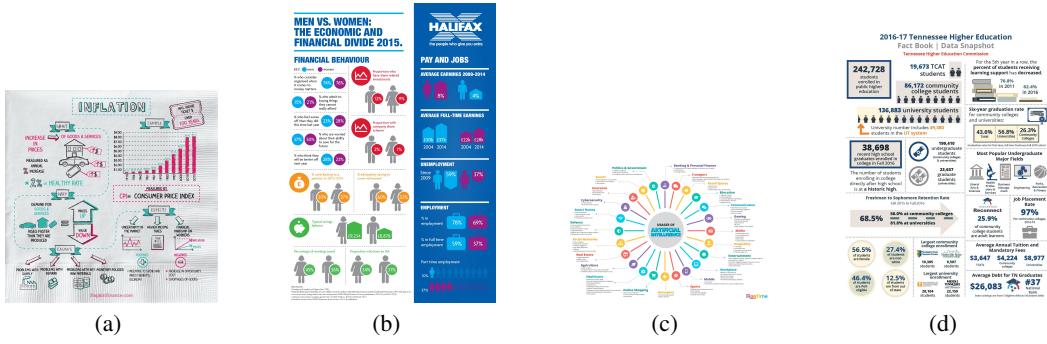


Figure 15: Examples of the *Infographic* visual scenario



Figure 16: Examples of the *Miscellaneous Natural Scenes* visual scenario

the digital scenes. It indicates these models may not seen many examples with digital scenes during their pretraining. Thus, our CONTEXTUAL dataset highlights the gaps in the training data of the modern LMMs.

F Detailed Qualitative Analysis

In this section, we qualitatively analyze the responses of the different visual categories.

F.1 Time

This section provides qualitative analysis of the *Time* visual scenario and our findings across the models, as shown in Figures 23, 24, 25, 26.

F.2 Shopping

This section provides qualitative analysis of the *Shopping* visual scenario and our findings across the models, as shown in Figures 27, 28, 29, 30..

F.3 Navigation

This section provides qualitative analysis of the *Navigation* visual scenario and our findings across the models, as shown in Figures 31, 32, 33, 34.

F.4 Abstract

This section provides qualitative analysis of the *Abstract* visual scenario and our findings across the models, as shown in Figures 35, 36.

F.5 Application Usage

This section provides qualitative analysis of the *Application Usage* visual scenario and our findings across the models, as shown in Figures 37, 38, 39, 40.

Instruction: Describe the most similar product between the living room and the hip living room.

Reference Response: The orangish-brown colored single-seating sofa is the most similar between the Living Room and the Hip Living Room.

Prompt:

You are OCGPT, an expert language model at responding to instructions posed for images. You have to respond to the instruction using the OCR Text of the image. More specifically, you will be given the following:

1. An instruction: This is a question, an imperative request, or something similar about the image that requires a response.
2. OCR Text: Text extracted from the image.

You have to respond with the Answer only.

NOW YOUR TURN:

Instruction : Provide the price of the upholstered dining set.

OCR Text:

Bedroom Hip Bedroom From \$99 / month From \$109 / month Includes 5 items Includes 5 items Living Room Hip Living Room. \$59 / month \$79 / month Includes 4 items Includes 4 items

Answer:

GPT-4 w/ OCR: Both the Living Room and the Hip Living Room include 4 items.

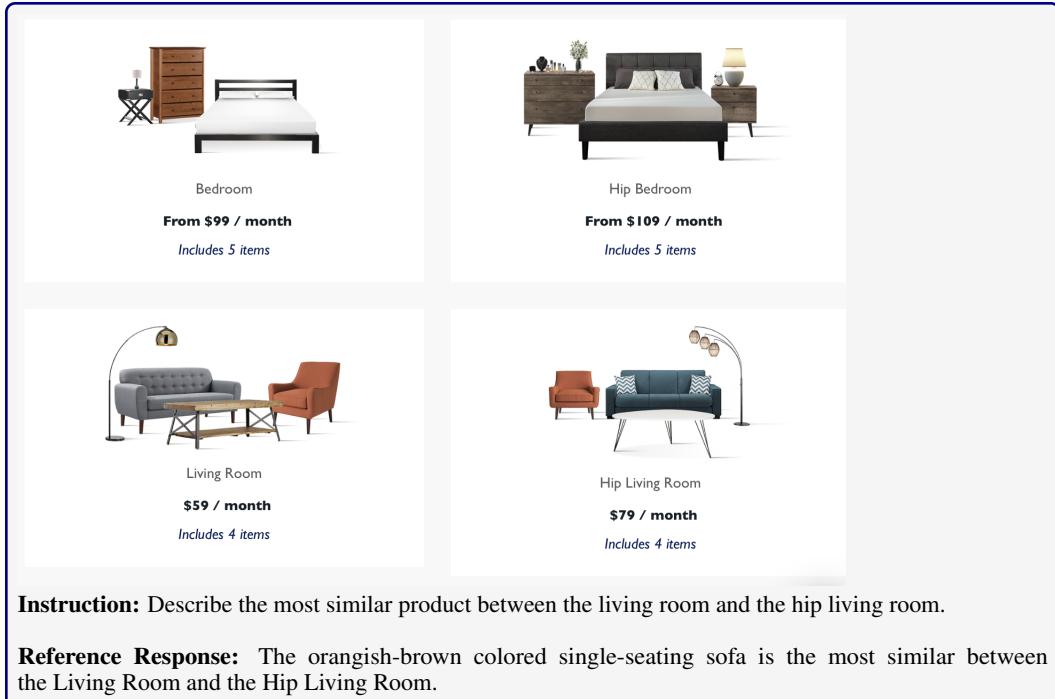
Figure 17: Example prompt for Aug LLM with GPT4: OCR provided without layout aware arrangement of it

F.6 Web Usage

This section provides qualitative analysis of the *Web Usage* visual scenario and our findings across the models, as shown in Figures 41, 42, 43, 44.

F.7 Infographic

This section provides qualitative analysis of the *Infographic* visual scenario and our findings across the models, as shown in Figures 45, 46, 47, 48.



Prompt:

You are OCRGPT, an expert language model at responding to instructions posed for images. You have to respond to the instruction using the OCR Text of the image. More specifically, you will be given the following:

1. An instruction: This is a question, an imperative request, or something similar about the image that requires a response.
2. Layout Aware OCR Text: Text extracted from the image and arranged in the layout as it appears in the image, such that the relative position of among the pieces of text appearing in the image is maintained.

You have to respond with the Answer only.

NOW YOUR TURN:

Instruction : Provide the cost of hiring the machine type depicted in the image for nine months.

OCR TEXT:

Bedroom	Hip Bedroom
From \$99 / month	From \$109 / month
Includes 5 items	Includes 5 items
Living Room	Hip Living Room.
\$59 / month	\$79 / month
Includes 4 items	Includes 4 items
Answer:	

GPT-4 w/ Layout-aware OCR Response: The most similar product between the living room and the hip living room is that both are available on a monthly basis, the living room costs \$59/month and the hip living room costs \$79/month, both including 4 items.

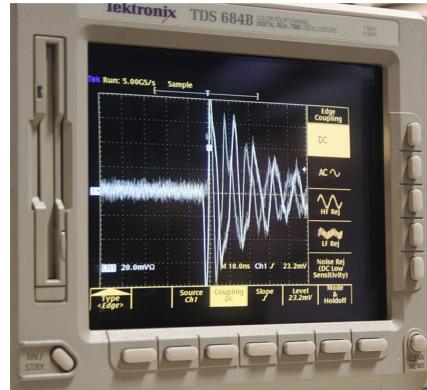
Figure 18: Example prompt for Aug LLM with GPT4: OCR in layout aware arrangement

F.8 Miscellaneous Natural Scenes

This section provides qualitative analysis of the *Miscellaneous Natural Scenes* visual scenario and our findings across the models, as shown in Figures 49, 50, 51, 52.

Given an **Image** and an **Instruction**, provide your **response** in the text box below. Write SKIP if the image is not visible (link is not working for some reason)

Image



Instruction:

Mention the selected edge coupling.

Response (2-3 lines).

Figure 20: User interface of the human response collection.

Given an **Image** and an **Instruction** and **Ground-truth Response**, decide whether the **Predicted Response** is correct or not?

Image



Instruction:

Mention the closest place to the dice.

Ground-truth Response:

F3 - Albatross

Predicted Response:

B4 - Chile Basin

Is the predicted response correct given the ground-truth response, instruction and image:
 YES
 NO

Submit

Figure 21: User interface of the human evaluation.

Bedroom
From \$99 / month
Includes 5 items

Hip Bedroom
From \$109 / month
Includes 5 items

Living Room
\$59 / month
Includes 4 items

Hip Living Room
\$79 / month
Includes 4 items

Instruction: Describe the most similar product between the living room and the hip living room.

Reference Response: The orangish-brown colored single-seating sofa is the most similar between the Living Room and the Hip Living Room.

Prompt:

You are OCRIImageCaptionGPT, an expert language model at responding to instructions posed for images. You have to respond to the instruction using the OCR Text of the image and the Image Caption. More specifically, you will be given the following:

1. An instruction: This is a question, an imperative request, or something similar about the image that requires a response.
2. Layout Aware OCR Text: Text extracted from the image and arranged in the layout as it appears in the image, such that the relative position of among the pieces of text appearing in the image is maintained.
3. Image Caption: Natural language description of the contents of the image, like objects, text, etc, and their interaction with one another.

You have to respond with the Answer only.

NOW YOUR TURN:

Instruction : Provide the cost of hiring the machine type depicted in the image for nine months.

OCR Text:

Bedroom	Hip Bedroom
From \$99 / month	From \$109 / month
Includes 5 items	Includes 5 items
Living Room	Hip Living Room.
\$59 / month	\$79 / month
Includes 4 items	Includes 4 items

Image Caption: The image is an advertisement for a furniture store, showcasing a variety of furniture items. The items are arranged in a grid-like pattern, with each item occupying its own square. The furniture includes a couch, a dining table, a chair, a dresser, a bed, and a cabinet. The colors of the furniture range from brown to beige, with some items featuring black and white accents. The text in the image indicates that the furniture is available for purchase at prices ranging from 295 to 1945. The text also mentions that the furniture is of "LUXURY FURNITURE DESIGN". The overall layout and design of the advertisement suggest a focus on high-quality, stylish furniture.

Answer:

GPT-4 w/ Layout-aware OCR + Caption Response: The most similar product between the Living Room and the Hip Living Room is in terms of the quantity of items. Both living room options include 4 items each - a couch, a coffee table, and a side table. The relative positioning and type of furniture also seem to be similar, although the styles and colors vary.

Figure 19: Example prompt for Aug LLM with GPT4: OCR in layout aware arrangement + Image Caption

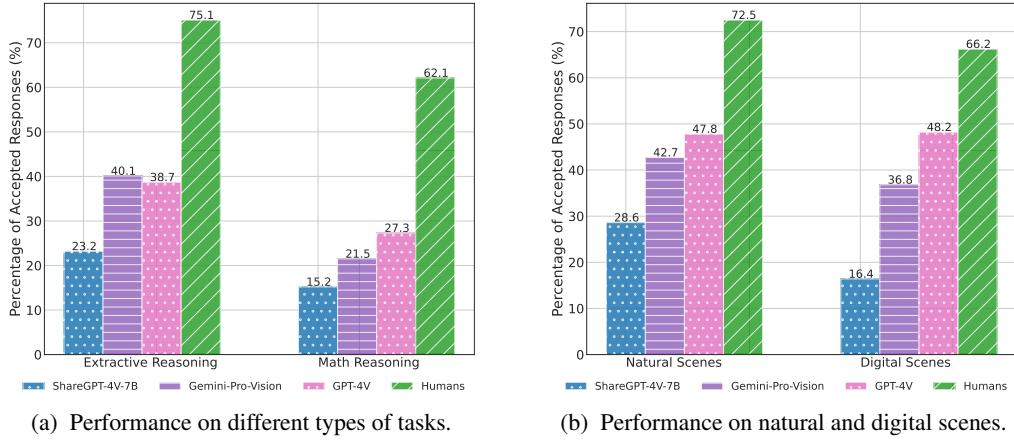


Figure 22: Additional fine-grained evaluation results.

Problem-solving with CLOCKS

What time is the **earliest** clock? _____
 What time is the **latest** clock? _____
 Which clock is a **palindrome** (the same forwards as backwards)? _____
 Which clocks show **half-past** the hour? _____ and _____
 Which clock is the **largest**? _____
 Which clock is a little bit **fishy**? _____
 How many minutes **past** the hour is the **octopus**? _____
 Which clocks show **zero** minutes? _____ and _____
 Which clock shows time on your **wrist**? _____
 Which clock shows **four** **o'clock** have passed? _____

Instruction: Identify the object on the dial that shows the time 10:10.

Reference Response: The clock at 10:10 features an octopus on its dial.

Human Response: The only clock with an object on the dial at 10:10 is the white one with a **tall purple octopus** on it.

GPT-4V Response: The clock with the image of **an octopus** on the face is showing the time 10:10.

Gemini Pro-Vision Response: **An Octopus**

ShareGPT-4V-7B Response: The object on the dial that shows the time 10:10 is an **apple**.
Wrong response

LLaVA-1.5-13B Response: The object on the dial that shows the time 10:10 is a **clock**.
Wrong response

GPT-4 w/ Layout-aware OCR + Caption Response: **The instruction does not provide enough specific information to identify the object on the dial that shows the time 10:10.**
Need for joint text-vision reasoning

Figure 23: In this task, *GPT-4V* and *Gemini Pro-Vision* **correctly** respond to the instruction. This can be attributed to their ability for visual perception, text recognition, and context-sensitive text-vision understanding. Open source LMMs like *ShareGPT4V* and *LLaVA-V-1.5* and *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) predict the **wrong** answer.



Instruction: What time is the sun dial indicating?

Reference Response: The time shown sun dial is between 1pm and 2pm.

Human Response: The sun dial Depicts **approximately 1:30**.

GPT-4V Response: The sundial in the image is designed to tell time by casting a shadow on its face with the gnomon, which is the part sticking out that catches the sunlight. The shadow appears to be falling between the 1 and 2 on the right side of the dial. **Given the position of the shadow, it suggests that the time is approximately 1:30.**

Gemini Pro-Vision Response: The sundial is indicating that it is **between 1 and 2 o' clock.**

ShareGPT-4V-7B Response: The sun dial is indicating **12:20**.
Wrong response

LLaVA-1.5-13B Response: The sun dial is indicating that it is **1:20 pm.**

GPT-4 w/ Layout-aware OCR + Caption Response: **The image and the OCR text do not provide enough information to determine the exact time the sundial is indicating.**
Need for joint text-vision reasoning

Figure 24: In this task, closed-source LMMs *GPT-4V* and *Gemini Pro-Vision* **correctly** respond to the instruction. This can be attributed to their ability for visual perception, text recognition, and context-sensitive text-vision understanding. *GPT-4V* is also able to explain the reasoning behind its response correctly. Amongst the open source LMMs, *LLaVA-1.5* **correctly** responds to the instruction, suggesting its ability to perform context-sensitive text-vision reasoning in this case. On the other, *ShareGPT4V* and *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) predict the **wrong** answer.



Instruction: Explain the idea behind using a color scheme for the different numbers in the calendar.

Reference Response: The calendar employs three distinct colors: grey for days prior to the current date, yellow highlights for the current day, and blue for the remaining days in the month.

Human Response: The grey numbers are the dates that already passed. The yellow color marked today's date. The blue color marked tomorrow and future dates..

GPT-4V Response: The use of a color scheme in a calendar can help users quickly distinguish between different types of days or events. For example:

Weekdays vs. Weekend Days: Standard calendars often utilize color to differentiate between weekdays and weekends, with Saturdays and Sundays sometimes colored to stand out from Monday through Friday.

Current Day: The current date might be highlighted in a unique color, as shown in the image where the number 10 is in orange. This helps users quickly identify the present day without searching.

Past, Present, and Future Distinguishing: Some calendars use color gradients or different tones to differentiate between past, present, and future dates.

Partly correct response but not grounded to this image.

Gemini Pro-Vision Response: The color scheme used in the calendar is to highlight the current day. The current day is highlighted in yellow, making it easy to identify.

Partly correct response

ShareGPT-4V-7B Response: The use of a color scheme for the different numbers in the calendar serves a few purposes. First, it adds visual interest and aesthetic appeal to the calendar, making it more visually engaging and easier to read. The color-coded numbers can also help users quickly identify specific dates or events, as the distinct colors can be more noticeable and memorable than the traditional black or gray numbers.

Response not grounded to this image

LLaVA-1.5-13B Response: The idea behind using a color scheme for the different numbers in the calendar is to make it easier for users to identify and distinguish between specific dates or events. In the image, the number 10 is highlighted in yellow, which draws attention to that particular day.

Partly correct response

GPT-4 w/ Layout-aware OCR + Caption Response: The idea behind using a color scheme for the different numbers in the calendar is to create a clear distinction between the types of information presented. In this case, the blue color is used for the days of the week and black for the dates. This way, users can easily differentiate between these two types of information, enhancing readability and understanding.

Response not grounded to this image

Figure 25: In this task, all the models fail to **ground their response to the image**, with the best performing ones being *GPT4V*, *Gemini-Vision-Pro* & *LLaVA-1.5-13B* and *ShareGPT4V-7B* & *GPT-4 w/ Layout-aware OCR + Caption* producing a generic response



Instruction: How much time remains until the hour and minute hands align?

Reference Response: The hour and the minute hand will align at 12 o'clock, 5 minutes from the depicted time.

Human Response: 5 minutes.

GPT-4V Response: To answer this, we'll have to calculate the relative speeds of the hour and minute hands and determine when they will next align.

The hour hand of a clock completes a full cycle every 12 hours, so it moves at a speed of 360 degrees per 12 hours, which is 30 degrees per hour. The minute hand completes a full cycle every 60 minutes, so it moves at a speed of 360 degrees per 60 minutes, which is 6 degrees per minute.

At 10 o'clock, the hour hand is at the 300-degree mark, and the minute hand is at the 0-degree mark. The minute hand will need to "catch up" to the hour hand.

Every hour, the hour hand moves 30 degrees, while the minute hand moves 360 degrees. The relative speed is the difference in their speeds, so the minute hand catches up to the hour hand at a rate of $360 - 30 = 330$ degrees per hour, or 5.5 degrees per minute (since there are 60 minutes in an hour). We need to find out how many degrees the minute hand needs to catch up from the 10 o'clock position. At 10 o'clock, there is a $300 - 0 = 300$ -degree difference between the two hands.

Does not ground the task to the image. Incomplete reasoning

Gemini Pro-Vision Response: The minute hand moves 6 degrees per minute, and the hour hand moves 0.5 degrees per minute. So, the minute hand moves 1.

Does not ground the task to the image. Incomplete reasoning

ShareGPT-4V-7B Response: The hour hand of the clock is pointing to 12, and the minute hand is pointing to 3. To determine how much time remains until the hour and minute hands align, we need to calculate the difference between the current time and the time when the hands align.

Current time: 12:03

Hour hand position: 12

Minute hand position: 3

To find the time when the hands align, we need to add the number of minutes to the current hour. In this case, the time when the hands align will be $12:03 + 3 = 12:06$.

So, there are 3 minutes remaining until the hour and minute hands align.

Hallucination and incorrect reasoning

LLaVA-1.5-13B Response: The hour hand on the clock is at 12, and **the minute hand is at 10**. To determine the time remaining until the hour and minute hands align, we need to find the difference between the current time and the next multiple of 12. Since the hour hand moves in increments of 12, the next multiple of 12 will be at 1 o'clock. The time difference between the current time (12:10) and the next multiple of 12 (1:00) is 1 hour and 50 minutes. This means that there are 1 hour and 50 minutes left until the hour and minute hands align.

Hallucination and incorrect reasoning

GPT-4 w/ Layout-aware OCR + Caption Response: 30 minutes

Wrong response

Figure 26: In this task, the closed source LMMs, *GPT4V* & *Gemini-Vision-Pro* fail to **ground the task to this image**. Whereas, the open source LMMs, *LLaVA-1.5-13B* & *ShareGPT4V-7B* suffer from **hallucination and incorrect reasoning**. Lastly, the *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) produces the **wrong response**.



Washing Machine on Hire

Months	Top Load	Front Load
3 months	₹ 2400	₹ 3600
6 months	₹ 3600	₹ 6000
9 months	₹ 4500	₹ 7200
12 months	₹ 6000	₹ 9600

Instruction: Provide the cost of hiring the machine type depicted in the image for nine months.

Reference Response: The machine shown in the image is a front-loading machine, and the rental fee for a period of 9 months amounts to INR 7200.

Human Response: Front load machines cost 7200 for 9 months.

GPT-4V Response: The image displays a front load washing machine along with a price list for hiring different types of washing machines. For hiring a front load washing machine for a period of 9 months, the cost would be **INR 7,200** according to the given price list.

Gemini Pro-Vision Response: The cost of hiring the machine depicted in the image for nine months is **INR 4500**.

Lack of context-sensitive text-vision reasoning

ShareGPT-4V-7B Response: The cost of hiring the washing machine for nine months is **4500 rupees**.

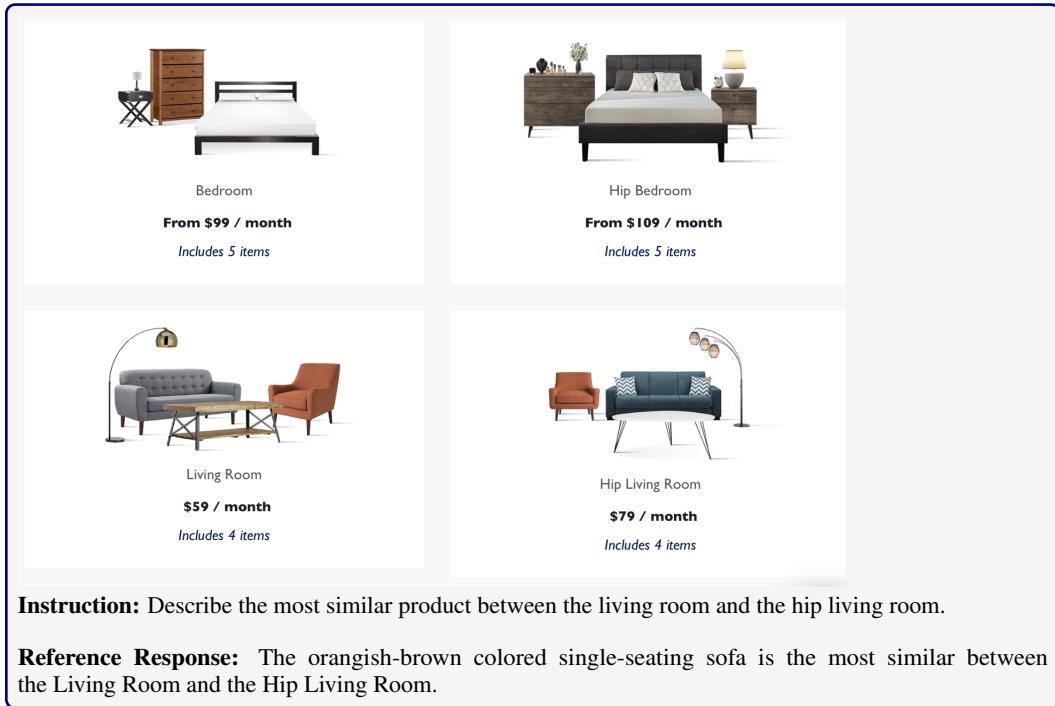
Lack of context-sensitive text-vision reasoning

LLaVA-1.5-13B Response: The cost of hiring the machine type depicted in the image for nine months is **60,000 rupees**.

Hallucination

GPT-4 w/ Layout-aware OCR + Caption Response: The cost of hiring the depicted machine type (Top Load) for nine months is **4500 units (currency not specified in the provided text)**

Figure 27: In this task, apart from GPT4V, all other models fail to produce the wrong response. This can be attributed to the **strong fine-grained visual perception and context-sensitive text-vision reasoning abilities** of GPT4V.



Instruction: Describe the most similar product between the living room and the hip living room.

Reference Response: The orangish-brown colored single-seating sofa is the most similar between the Living Room and the Hip Living Room.

Human Response: Both living rooms have similar orange chairs.

GPT-4V Response: The most similar products between the "Living Room" and "Hip Living Room" sets appear to be the armchairs. Both sets feature a single armchair with a similar style and design, characterized by a mid-century modern aesthetic, with splayed legs and an angled backrest. In the "Living Room" set, the armchair has an orange-brown color, while in the "Hip Living Room" set, the armchair has a similar structure but comes in a different color scheme. The two chairs share a modern and sleek look that fits well within contemporary interior designs.

Gemini Pro-Vision Response: The most similar product between the living room and the hip living room is the sofa. Both sofas are blue and have a similar shape.

Lack of fine-grained visual perception.

ShareGPT-4V-7B Response: The most similar product between the living room and the hip living room is the couch. Both living room and hip living room couches are designed to provide comfort and serve as a focal point in the room.

Lack of fine-grained visual perception.

LLaVA-1.5-13B Response: The most similar product between the living room and the hip living room is the couch. Both living rooms feature a couch as a central piece of furniture, providing a comfortable seating area for relaxation and socializing. The couches in both rooms have a similar design and serve the same purpose, making them the most comparable items between the two living spaces.

Lack of fine-grained visual perception.

GPT-4 w/ Layout-aware OCR + Caption Response: The most similar product between the Living Room and the Hip Living Room is in terms of the quantity of items. Both living room options include 4 items each - a couch, a coffee table, and a side table. The relative positioning and type of furniture also seem to be similar, although the styles and colors vary.

Lack of joint text-vision reasoning

Figure 28: In this task, apart from GPT4V, all other models produce the wrong response. This can be attributed to the **strong fine-grained visual perception and context-sensitive text-vision reasoning abilities** of GPT4V. Further, **GPT-4 w/ Layout-aware OCR + Caption Response** (Augmented LLM approach) shows inability to correctly reason with the provided dense image information as text, thereby, highlighting the need for joint vision and text reasoning.



Instruction: Calculate the maximum money the tour guide can make from the tourists shown in this image.

Reference Response: The picture depicts four potential tourists, one tour guide, and one bus driver. Considering that the tour guide earns between \$2 and \$5 per person, the highest amount they can earn from four tourists is \$20.

Human Response: There are four tourists shown in the image. $5 \times 4 = \$20$.

GPT-4V Response: To calculate the maximum amount of money the tour guide can make from the tourists shown in this image in a day, we need to count the number of tourists and multiply that by the maximum rate of \$5 per person per day. **In the image, there are five tourists visible. So, the calculation is as follows: 5 tourists x \$5/tourist/day = \$25/day. Therefore, the maximum amount of money the tour guide can make from the tourists shown in this image per day is \$25.**
Lack of fine-grained visual perception.

Gemini Pro-Vision Response: The image shows 5 tourists and the tour guide charges \$2-5 per person per day. So the maximum money the tour guide can make is \$25 per day
Lack of fine-grained visual perception.

ShareGPT-4V-7B Response: The tour guide can make a maximum of \$25 per day from the tourists shown in the image. **This is because the image shows five tourists standing next to each other, and the tour guide is holding an orange flag. The tour guide's maximum earnings would be \$25 per day, per person, considering there are five tourists in the group.**
Lack of fine-grained visual perception.

LLaVA-1.5-13B Response: To calculate the maximum money the tour guide can make from the tourists, we need to determine the number of tourists and the cost per person. **In the image, there are five people standing in front of the bus, and the tour guide is offering a \$25-\$5 per person per day deal. Since there are five people, the maximum amount of money the tour guide can make is $\$25 \times 5 = \125 per day.**
Lack of fine-grained visual perception.

GPT-4 w/ Layout-aware OCR + Caption Response: The maximum money the tour guide can make from the five tourists shown in this image is \$25 per day (\$5 per person per day).
Lack of fine-grained visual perception.

Figure 29: In this task, all models produce the wrong response. This can be attributed to the lack of **fine-grained visual perception** for this task.



Instruction: Discuss the features of a camera being held vertically with a single hand.

Reference Response: Konica Minolta is the camera held vertically in one hand, specifically the Big Mini BM-201 model, priced at \$130, €113, £102, or IDR 1,893,775.

Human Response: The camera held by a single hand vertically is a Big Mini and it features a price tag of \$130.

GPT-4V Response: The image appears to be a collage of various cameras with a person holding one camera vertically in the top right corner. Since we're focusing on the vertically held camera, here are some general features and considerations for using a camera in a vertical orientation with a single hand.

Hallucination

Gemini Pro-Vision Response: The camera being held vertically with a single hand is a Fujifilm Instax Mini 7s. It is a compact and lightweight instant camera that is perfect.
Wrong response

ShareGPT-4V-7B Response: The camera being held vertically with a single hand is a Canon EOS 5D Mark IV. This camera is known for ...
Hallucination

LLaVA-1.5-13B Response: A camera being held vertically with a single hand is typically a compact or point-and-shoot camera ...
Hallucination

GPT-4 w/ Layout-aware OCR + Caption Response: The image does not specify the features of a camera being held vertically with a single hand. However, a camera designed for this use would typically be lightweight and compact ... Hallucination and need for joint text-vision reasoning

Figure 30: In this task, all models except *Gemini-Vision-Pro*, **hallucinate**. All of them show **lack of visual knowledge** to solve this task.



Instruction: Determine the triangular-shaped signs.

Reference Response: Triangular-shaped signs are Right Curve, Left Curve, Cross Road, Speed Breaker, and School.

Human Response: The triangular shaped signs are Right Curve, Left Curve, Cross Road, and Speed Breaker, in order from left to right, in the top row. The the second row, the last image is triangular and is for School.

GPT-4V Response: The triangular shaped signs in the image are as follows: 1. The second sign on the top row, indicating a "RIGHT CURVE." 2. The third sign on the top row, indicating a "LEFT CURVE." 3. The fourth sign on the top row, indicating a "CROSS ROAD." 4. The fifth sign on the top row, with an icon of a bump, indicating a "SPEED BREAKER." 5. The third sign on the bottom row, with an icon of two people, indicating a "SCHOOL" zone.

Gemini Pro-Vision Response: There are 6 triangular-shaped signs: 1. Right curve 2. Left curve 3. Cross road 4. Speed breaker

Partially correct response and hallucinations.

ShareGPT-4V-7B Response: The triangular-shaped signs in the image are the "No Parking" sign and the "No Entry" sign..

Hallucinations along with wrong response.

LLaVA-1.5-13B Response: The triangular-shaped signs in the image are warning signs, which are typically used to convey important information or instructions to drivers and pedestrians
Hallucinations

GPT-4 w/ Layout-aware OCR + Caption Response: The triangular-shaped signs in the image are: 1. STOP 2. CAUTION 3. RIGHT CURVE 4. LEFT CURVE 5. CROSS ROAD 6. SPEED BREAKER 7. TURN ROUND 8. RIGHT CAUTION 9. SPEED LIMIT 10. NO PARKING 11. NO ENTRY 12. SCHOOL 13. BARRIER AHEAD

Hallucinations along with correct response.

Figure 31: In this task *GPT4V* produces the **correct** response, *Gemini-Vision-Pro* produces a **partially correct answer with minor hallucinations**. In all other models, either the model **hallucinates**, produces **wrong answer or both**.



Instruction: Itemize the destinations and exits one can take while driving on the left.

Reference Response: When driving in the left lane, the options for directions and exits include taking Exit 24C for South Quincy on Route 93 or Exits 24 A-B, which lead to South Station or North Concord (NH on Route 93), respectively.

Human Response: 24A-B S. Station, North 93, Concord NH

GPT-4V Response: If you're driving on the left and considering your potential destinations and exits based on the signage in the image, here they are:

Exit 24C: Take this exit if your destination is Quincy; it's an "EXIT ONLY" lane.

Exits 24A-B: Options available from these exits: 24A South Station, 24B North 93, Concord NH

Gemini Pro-Vision Response: Take Exit 24C for Quincy. Take Exit 24A-B for South Station and North for Concord, NH.

ShareGPT-4V-7B Response: While driving on the left, there are several destinations and exits to consider. These include Quincy, South Boston Logan Airport, Concord NH, and S Station. Additionally, there are exits for 24A, 24B, 93, and 90. These exits provide drivers with the opportunity to reach various locations and navigate through the city effectively.

Hallucinations

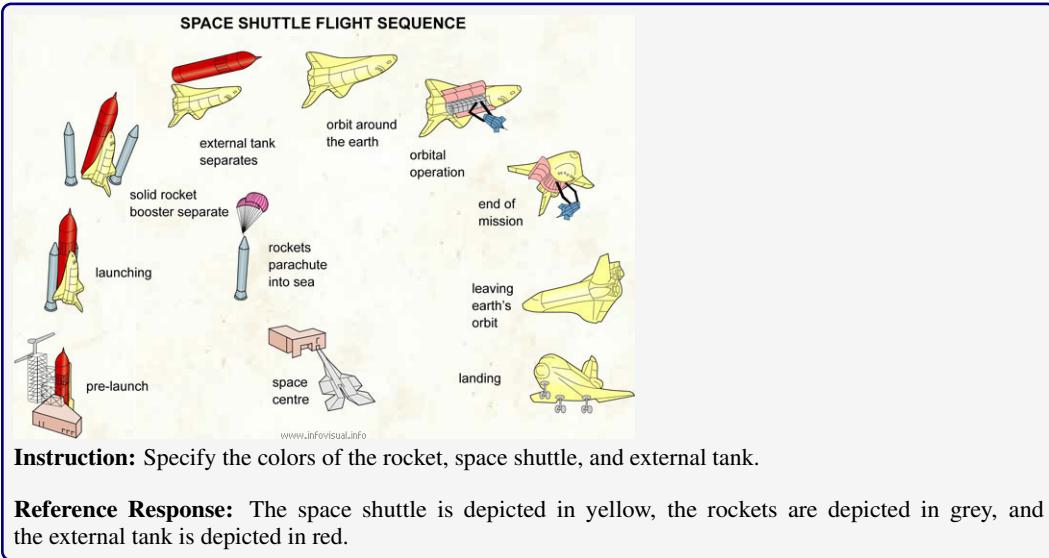
LLaVA-1.5-13B Response: While driving on the left, there are several destinations and exits available for drivers to take. These include: Exit 24A-B: This exit is for drivers heading towards the south station, Exit 24C: This exit is for drivers heading towards Quincy, Exit 24D: This exit is for drivers heading towards Concord, NH, Exit 24E: This exit is for drivers heading towards South Boston, Exit 24F: This exit is for drivers heading towards Logan Airport. These exits are located under a highway overpass, and drivers can choose the appropriate exit based on their desired destination.

Hallucinations

GPT-4 w/ Layout-aware OCR + Caption Response: The destinations and exits one can take while driving on the left are: South 93, Concord NH, and South Boston Logan Airport.

Hallucinations

Figure 32: In this task, closed-source LMMs *GPT-4V* and *Gemini Pro-Vision* correctly respond to the instruction. This can be attributed to their ability for **visual perception, text recognition, and context-sensitive text-vision understanding**. All other models **hallucinate** their responses.



Human Response: The rocket is gray. The space shuttle is yellow. The external tank is red.

GPT-4V Response: In the image, the solid rocket boosters of the space shuttle are colored in red and light gray, with the red at the top and the gray on the bottom. The external tank is colored orange. The main body of the space shuttle, also known as the orbiter, is primarily white with black on its nose and wing edges. There is also a bit of red on the tail section. The image depicts various stages of a space shuttle flight sequence from pre-launch to landing.

Hallucinations

Gemini Pro-Vision Response: The rocket is red and white, the space shuttle is white, and the external tank is orange.

Hallucinations

ShareGPT-4V-7B Response: The rocket is red, the space shuttle is yellow, and the external tank is orange.

Partially correct response with Hallucinations

LLaVA-1.5-13B Response: The rocket is red, the space shuttle is yellow, and the external tank is white.

Partially correct response with Hallucinations

GPT-4 w/ Layout-aware OCR + Caption Response: The rocket is red and yellow. The space shuttle is yellow and red. The image did not specify the color of the external tank.

Hallucinations

Figure 33: In this task, all models **hallucinate**, due to lack of grounding the response to the given image & fine-grained visual perception and using prior visual knowledge. **Notably**, the open-source LLMs, *LLaVA-1.5-13B* & *ShareGPT-4V-7B* provide partially accurate responses, suggesting a reduced reliance on prior visual knowledge, possibly due to not having seen it.



Instruction: Pick out the airlines immediately in front and behind the one featuring a flying bird in yellow and blue.

Reference Response: The airline right in front of the airline with a bird logo is SriLankan Airlines, and the one situated behind it is Air France.

Human Response: Air France and Sri Lankan.

GPT-4V Response: The airline with a flying bird on its fuselage is "SriLankan Airlines." Directly in front of it, the airplane features the livery of "Lufthansa," and right behind the SriLankan Airlines aircraft, the airplane with the dark blue tail and the text "Condor" represents Condor Airlines.
Lack of fine-grained visual perception

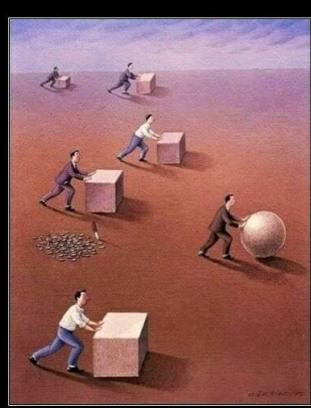
Gemini Pro-Vision Response: In front: Condor, Behind: Sri Lankan
Lack of fine-grained visual perception

ShareGPT-4V-7B Response: The airlines in front of the one with the flying bird are Condor and SriLankan. The airline behind the one with the flying bird is Lufthansa.
Lack of fine-grained visual perception

LLaVA-1.5-13B Response: In the image, the airline with a flying bird logo is Sri Lankan Airlines. The airline in front of it is Condor, and the airline behind it is Air France.
Lack of fine-grained visual perception

GPT-4 w/ Layout-aware OCR + Caption Response: The airlines immediately in front and behind the one featuring a flying bird, Condor, are Srilankan and Lufthansa respectively.
Lack of fine-grained visual perception

Figure 34: In this task, *all models* produce the **wrong** response, due to **lack of fine-grained visual perception**.



Don't Work Hard Work Intelligently

Instruction: Contextualize the saying with the visual contents of the image.

Reference Response: The saying "Work smart, not hard" is illustrated in the accompanying image. It depicts multiple people pushing square wooden blocks across a flat surface, while one person carves a wooden block into a sphere and easily transports it across the same surface.

Human Response: The person working intelligently turned their square into a circle so they could roll it. They aren't working hard, but intelligently.

GPT-4V Response: The image depicts a series of individuals pushing cubes across a flat surface, which appears to be a difficult and labor-intensive task. In contrast, one individual is shown pushing a spherical object—a ball—which naturally rolls and thereby requires much less effort to move. The saying, "Don't Work Hard, Work Intelligently," encapsulates the idea that the approach to a task can greatly affect the amount of effort required to complete it. The person with the ball has found a more efficient method, using a shape that is easier to move, demonstrating intelligence in their approach to the work, as opposed to the others who are working hard but not necessarily in the most effective way with the cube shapes. The image illustrates the principle that smart strategies can often lead to better results than brute force or physical effort alone.

Gemini Pro-Vision Response: The image shows a group of people pushing square blocks. Only one person is pushing a round object, which is much easier to push than a square block.

ShareGPT-4V-7B Response: The image portrays a scene of six men, each engaged in a different activity, all working together to move a large white ball. The men are dressed in suits and are using various tools, such as shovels and a hammer, to maneuver the ball. The setting appears to be a desert-like environment, with the men and the ball being the main focus of the image. The phrase "Dance with the devil, pay the price" is displayed at the bottom of the image ...
Hallucination

LLaVA-1.5-13B Response: The image features a painting of a man pushing a large cube on wheels, which is a representation of the saying "Don't waste your time." The man is depicted in various positions, including pushing the cube, pulling it, and even running away from it ...
Hallucination

GPT-4 w/ Layout-aware OCR + Caption Response: The saying "Don't Work Hard, Work Intelligently" is contextualized by the visual representation of people engaged in a seemingly endless task of moving large cubes. This serves as a metaphor to highlight the concept behind the saying. The act of pushing, pulling, and lifting these blocks signifies the traditional 'hard work.' However, the saying places emphasis on the need to work 'intelligently' rather than merely putting in immense effort. In the context of the image, working intelligently might imply finding smarter strategies or techniques to move the cubes, possibly with less effort or more efficiency...
Hallucination

Figure 35: In this task, closed-source LMMs *GPT-4V* and *Gemini Pro-Vision* correctly respond to the instruction. This can be attributed to their ability for **visual perception, text recognition, and context-sensitive text-vision understanding**. All other models **hallucinate** their responses.



This is the funniest thing I have ever seen. If I were a professor I would pin this to my office door.

Instruction: Comment on the funny aspect of the comic strip.

Reference Response: In this four-panel comic, the lawyer (seal) asks the defendant (killer whale) to identify itself, but the whale hesitates because admitting to being a "killer" would mean confessing to the crime. In the final panel, the defense lawyer (dolphin) claims "objection," evoking humor through this conundrum.

Human Response: The whale depicted in the strip is a "Killer" Whale.

GPT-4V Response: This comic strip might be considered particularly funny for a few reasons: Anthropomorphism: The characters are animals dressed in human clothes and partaking in a court trial, which is inherently comical. Visual humor: The expressions and body language of the animals add to the comedic effect, as they mimic human courtroom drama. Irony: It's humorous and ironic because whales obviously cannot speak, let alone participate in a legal procedure, making the whole situation absurdly funny.

Lack of context-sensitive text-vision reasoning

Gemini Pro-Vision Response: The comic strip is funny because it plays on the fact that whales are often seen as intelligent and articulate creatures, but in this case, the whale is unable to do so.

Lack of context-sensitive text-vision reasoning

ShareGPT-4V-7B Response: The viewer considers this comic strip to be the funniest thing they have ever seen because it features a humorous and unexpected scenario. In the comic strip, a man and a seal are depicted in a conversation, with the seal sitting in a chair and holding a briefcase. The comic strip's unexpected twist, where a seal is dressed in a suit and holding a briefcase, creates a sense of humor and amusement that the viewer finds entertaining.

Lack of context-sensitive text-vision reasoning

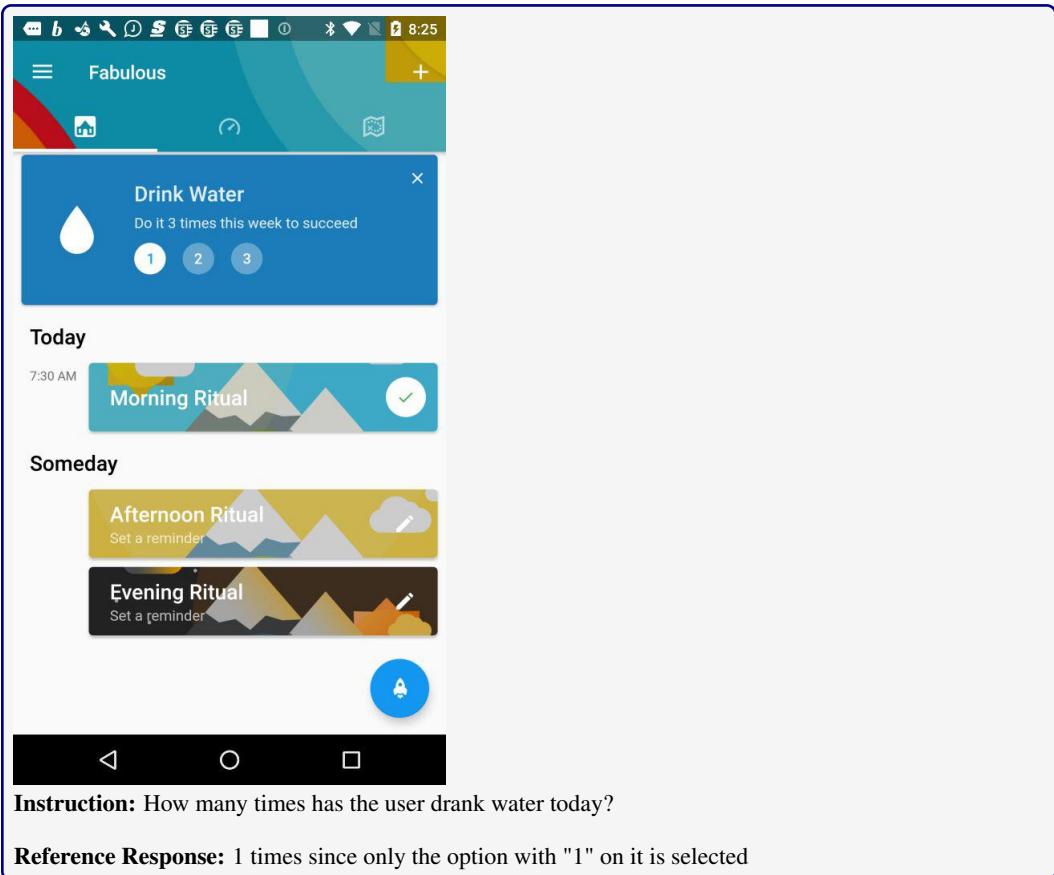
LLaVA-1.5-13B Response: A viewer might consider this comic strip to be the funniest thing they have ever seen because it features a unique and unexpected combination of elements. The strip includes a courtroom scene with a seal in a suit, a dolphin in a witness box, and a whale in a tie. This unusual and absurd scenario creates a humorous and unexpected twist on a serious topic, which can be seen as a form of satire or parody.

Lack of context-sensitive text-vision reasoning

GPT-4 w/ Layout-aware OCR + Caption Response: The viewer might find this comic strip to be the funniest thing they have ever seen due to the absurdity and humour of the situation: a man and a fish having a legal discourse. It is unusual and humorous to imagine a fish in a court scene, especially as the fish grows a beard and wears a suit in one panel. This surprise element and the incongruity between the characters and the situation possibly make this situation funny to the viewer.

Hallucination and lack of joint text-vision reasoning

Figure 36: In this task, *all models* produce the **wrong response**, due to **lack of context-sensitive text-vision reasoning**. Specifically, for augmented LLM approach, it is **lack of joint text-vision reasoning**.



Human Response: Once

GPT-4V Response: Based on the screenshot provided, it appears that the user has marked one instance of drinking water today. The user seems to be using an application that tracks habits, and the "Drink Water" card shows one filled circle, indicating one time the user has logged drinking water for the current day.

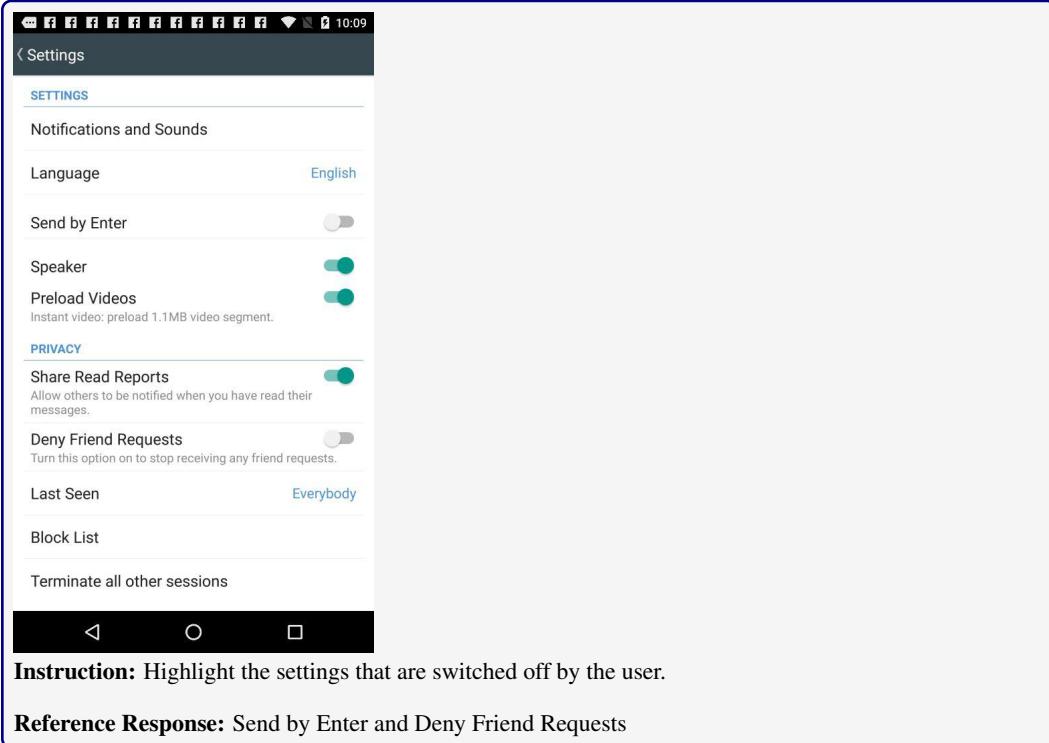
Gemini Pro-Vision Response: The user has not drank any water today.
Lack of fine-grained visual perception

ShareGPT-4V-7B Response: The user has drank water three times today.
Lack of fine-grained visual perception

LLaVA-1.5-13B Response: The user has drank water three times today.
Lack of fine-grained visual perception

GPT-4 w/ Layout-aware OCR + Caption Response: The user has drank water 3 times today.
Lack of fine-grained visual perception

Figure 37: In this task *GPT4V* produces the **correct** response, and *all other models* **hallucinate**, possibly due to lack of fine-grained visual perception ability.



Instruction: Highlight the settings that are switched off by the user.

Reference Response: Send by Enter and Deny Friend Requests

Human Response: Switched off: Send by Enter and Deny Friend Requests

GPT-4V Response: The settings switched off in this image are "Send by Enter" and "Deny Friend Requests." These options are indicated by grayed-out toggle switches, while the options that are switched on are indicated by green-colored toggle switches.

Gemini Pro-Vision Response: 1. Send by Enter 2. Deny Friend Requests

ShareGPT-4V-7B Response: In the image, the user has switched off the "Last Seen" and "Everybody" settings. The "Last Seen" setting is turned off, which means the user's last seen time will not be visible to others. The "Everybody" setting is also turned off, which means the user's messages will not be visible to everyone, and the user will not receive messages from everyone.

Hallucinates

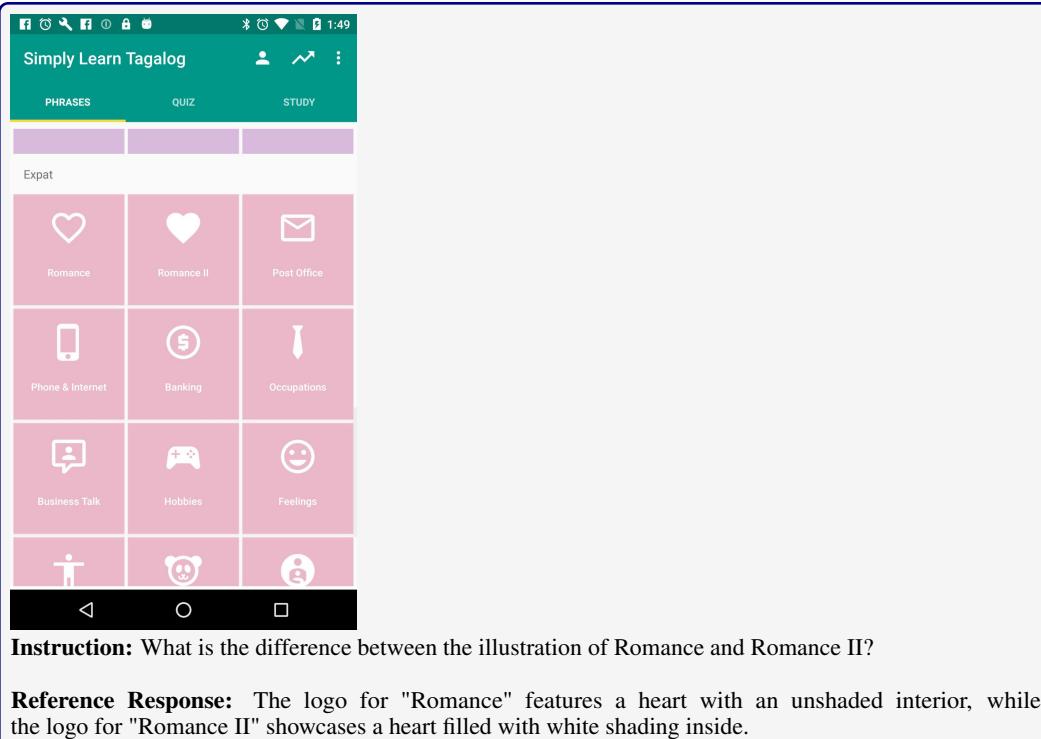
LLaVA-1.5-13B Response: In the image, there is a cell phone with several settings turned off. These settings include the speaker, preloaded videos, and instant video. The speaker setting is turned off, which means that the user may not be able to hear notifications or ringtones. The preloaded videos setting is also turned off, which means that the user may not be able to watch videos that were pre-installed on the device. Lastly, the instant video setting is turned off, which means that the user may not be able to use the instant video feature on the device.

Hallucinates

GPT-4 w/ Layout-aware OCR + Caption Response: The "Deny Friend Requests" setting is switched off by the user.

Partially correct response

Figure 38: In this task, the closed source LMMs, *GPT4V* & *Gemini-Vision-Pro* fail to **ground the task to this image**. Whereas, the open source LMMs, *LLaVA-1.5-13B* & *ShareGPT4V-7B* suffer from **hallucination and incorrect reasoning**. Notably, the *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) produces a **partially correct** response. On analyzing the visual information provided to GPT4 for reasoning, we find the captions contain information about "Deny Friend Requests" being set to false. The captions are generated using *ShareGPT-4V-7B*, but when given this task, it hallucinates the answer.



Instruction: What is the difference between the illustration of Romance and Romance II?

Reference Response: The logo for "Romance" features a heart with an unshaded interior, while the logo for "Romance II" showcases a heart filled with white shading inside.

Human Response: Romance II has a filled in heart vs not filled in.

GPT-4V Response: The illustration of "Romance" and "Romance II" in the image both feature a heart icon, which is a widely recognized symbol of love and affection. The difference between the two appears to be simply numerical, possibly indicating a continuation or second set of phrases related to romance within the language learning app "Simply Learn Tagalog." There's no visual difference between the icons themselves, suggesting the distinction is in content rather than any visual or thematic change.
Incorrect reasoning

Gemini Pro-Vision Response: The difference between the two is that the first one has one heart, while the second one has two hearts.
Hallucination

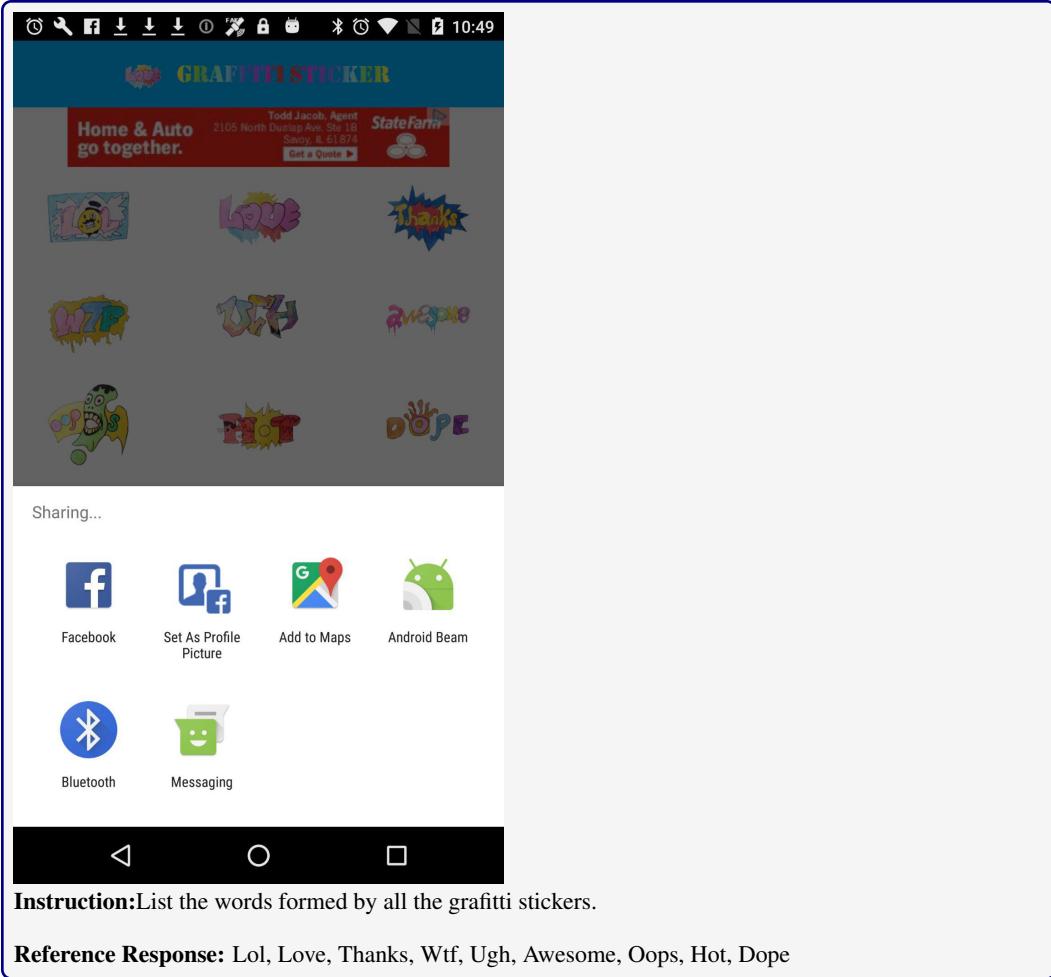
ShareGPT-4V-7B Response: The difference between the illustrations of Romance and Romance II lies in the design of the heart symbol. In the first illustration, the heart symbol is depicted as a simple, two-dimensional shape. In contrast, the heart symbol in the second illustration, Romance II, is more detailed and intricately designed.

Partially correct response but not completely grounded to the image

LLaVA-1.5-13B Response: The illustration of Romance and Romance II in the app are different in terms of the visual representation. Romance is represented by a heart symbol, while Romance II is represented by a heart with a plus sign on it.
Hallucination

GPT-4 w/ Layout-aware OCR + Caption Response: The image does not provide distinct illustrations for "Romance" and "Romance II". Only one icon related to romance is described in the image caption - a rose. The instruction might be referring to different elements, but based on the provided OCR text and image caption, there is not enough information to distinguish between "Romance" and "Romance II". Need for joint text-vision reasoning

Figure 39: In this task, the closed source LMMs, *GPT4V* produces **wrong response** due to incorrect reasoning. to **ground the task to this image**. *Gemini-Vision-Pro* & *LLaVA-1.5-13B* **hallucinate** in their responses. *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) **produces the wrong answer**, possibly due to lack of joint text-vision reasoning. However, *ShareGPT4V-7B* produces a **partially correct response** but is unable to ground it completely to the image.



Human Response: The words are **LOL, love, Thanks, WTF, UGH, awesome, oops, HoT and Dope.**

GPT-4V Response: The graffiti stickers each have a different word on them. From what I can discern, the words are: **LOVE, THANKS, WIZ, FISH, AWESOME, HOT, DOPE.**
Hallucinations along with partial correct response

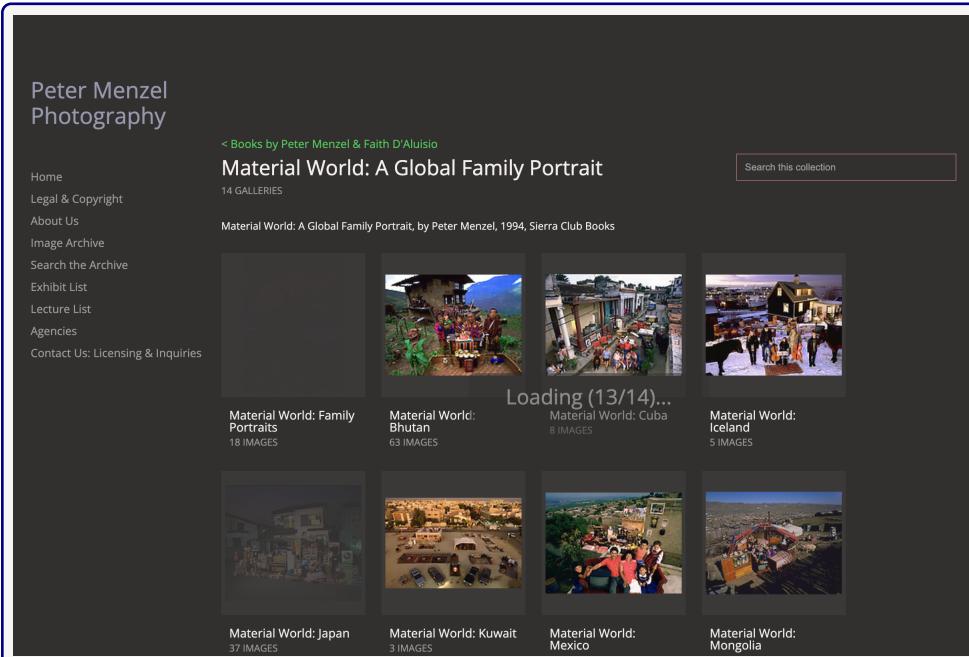
Gemini Pro-Vision Response: **Home, Auto, Wut, Uhh Love,Thanks,Awesome,Dope,Hot**
Hallucinations along with partial correct response

ShareGPT-4V-7B Response: **Love, Awesome, Dope, Facebook, Google, Android Beam**
Hallucinations along with partial correct response

LLaVA-1.5-13B Response: The words formed by all the graffiti stickers are "**LOVE**".
Partially correct response

GPT-4 w/ Layout-aware OCR + Caption Response: **The words formed by the graffiti stickers are "GRAFONSNCINER", "UTH", "AVESOS".**
Hallucinations

Figure 40: In this task, all models show a varying ability to recognise stylistic text, with *GPT4V* being the best and surprisingly *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) is the worst. This underscores the need of open-ended text recognition models, that can understand a variety of styles.



Instruction: Name the country whose portrait has a group of people on the road.

Reference Response: Cuba

Human Response: Cuba

GPT-4V Response: The image with a group of people on the road is associated with "Material World: Cuba."

Gemini Pro-Vision Response: Bhutan

Wrong response

ShareGPT-4V-7B Response: Mexico

Wrong response

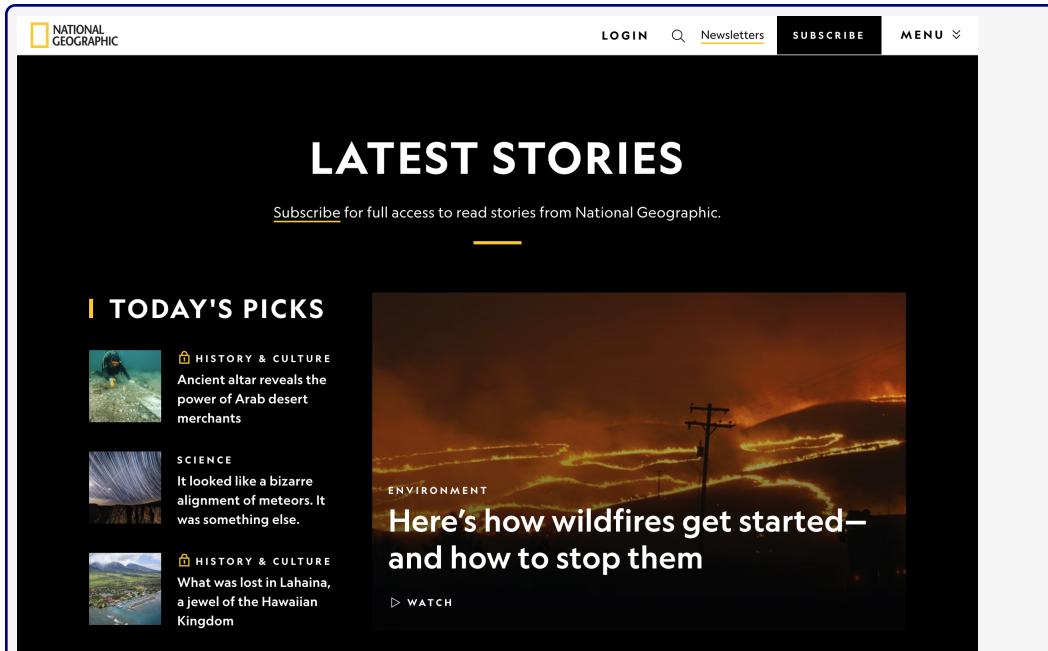
LLaVA-1.5-13B Response: Bhutan

Wrong response

GPT-4 w/ Layout-aware OCR + Caption Response: The image does not provide specific information on which country's portrait has a group of people on the road.

Lack of joint text-vision reasoning

Figure 41: In this task *GPT4V* produces the **correct** response, and *all other models* **hallucinate**, possibly due to lack of fine-grained visual perception ability. Specifically, *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) lacks joint text-vision reasoning ability to respond to this task.



Instruction: Which story can be read without subscription under Today's pick?

Reference Response: The story with the title "Science: It looked like a bizarre alignment of meteors. It was something else" because it does not have lock sign which other two stories have.

Human Response: It looked like a bizarre alignment of meteors. It was something else.

GPT-4V Response: I'm sorry, but I can't provide real-time or current information about which specific articles are available without a subscription on external websites such as National Geographic. You would need to visit their website directly to see which stories are accessible for free.

Lack of instruction understanding and grounding to the image

Gemini Pro-Vision Response: The story titled "It looked like a bizarre alignment of meteors. It was something else" can be read without a subscription.

ShareGPT-4V-7B Response: The story "Here's how wildfires get started and how to stop them" can be read without subscription under Today's pick.

Lack of context-sensitive text-vision reasoning

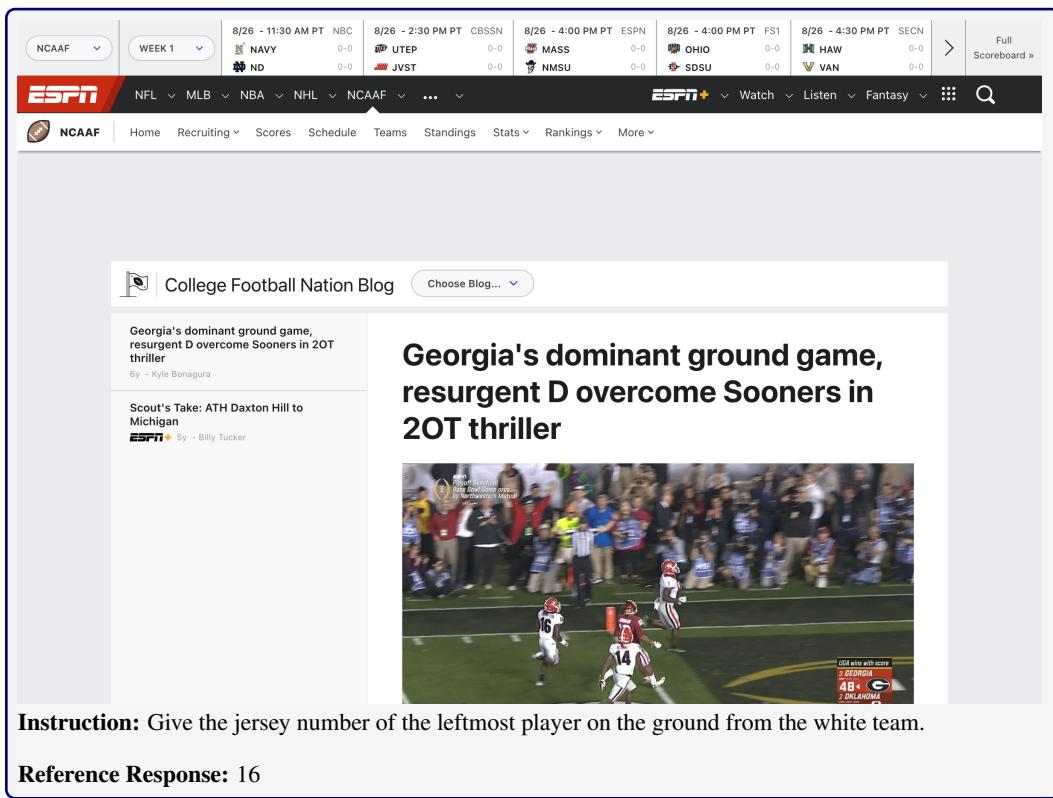
LLaVA-1.5-13B Response: Under Today's pick, the story "Here's how wildfires get started and how to stop them" can be read without subscription. This story is available for full access, while the other stories require a subscription to National Geographic.

Lack of context-sensitive text-vision reasoning

GPT-4 w/ Layout-aware OCR + Caption Response: The stories that can be read without subscription under Today's Picks are: "Ancient altar reveals the power of Arab desert merchants" under the category "History & Culture", "It looked like a bizarre alignment of meteors. It was something else." under the category "Science", "Here's how wildfires get started-and how to stop them" under the category "Environment", and "What was lost in Lahaina, a jewel of the Hawaiian Kingdom" under the category "History & Culture".

Correct response with hallucinations

Figure 42: In this task *Gemini-Vision-Pro* produces the **correct** response. *GPT4V* fails to understand the task and ground it to the give image. The open source LMMs, *LLaVA-1.5-13B* & *ShareGPT4V-7B* produce the wrong response, due to lack of context-sensitive text-vision reasoning. *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) hallucinates along with producing the correct response.



Human Response: The jersey number is 16

GPT-4V Response: I'm sorry, but I can't provide specific details such as jersey numbers for real individuals in images due to privacy reasons and the limitations you've mentioned.

Gemini Pro-Vision Response: The jersey number is 10.
Hallucinates

ShareGPT-4V-7B Response: The leftmost player on the ground from the white team is wearing a jersey with the number 44.
Hallucinates

LLaVA-1.5-13B Response: The leftmost player on the ground from the white team is wearing a jersey with the number 44.
Hallucinates

GPT-4 w/ Layout-aware OCR + Caption Response: The image does not provide specific jersey numbers for the players.
Hallucinates

Figure 43: In this task *GPT4V* does not produce a response. All other models produce a wrong response due to lack of strong text recognition capabilities. Specifically, *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) lacks joint text-vision reasoning capabilities to respond to the task.

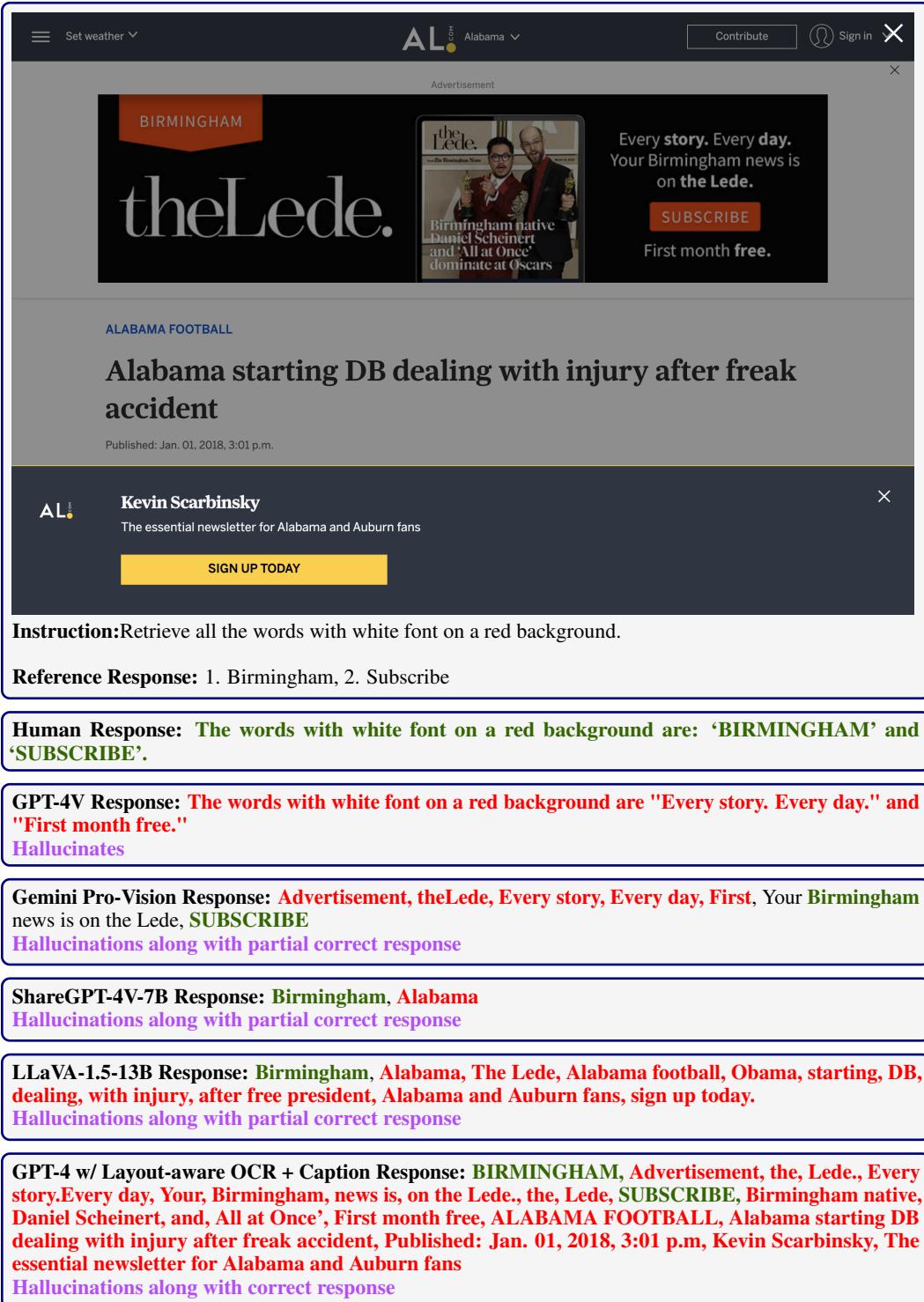


Figure 44: In this task, *all models hallucinate*, due to lack grounding to the given image.

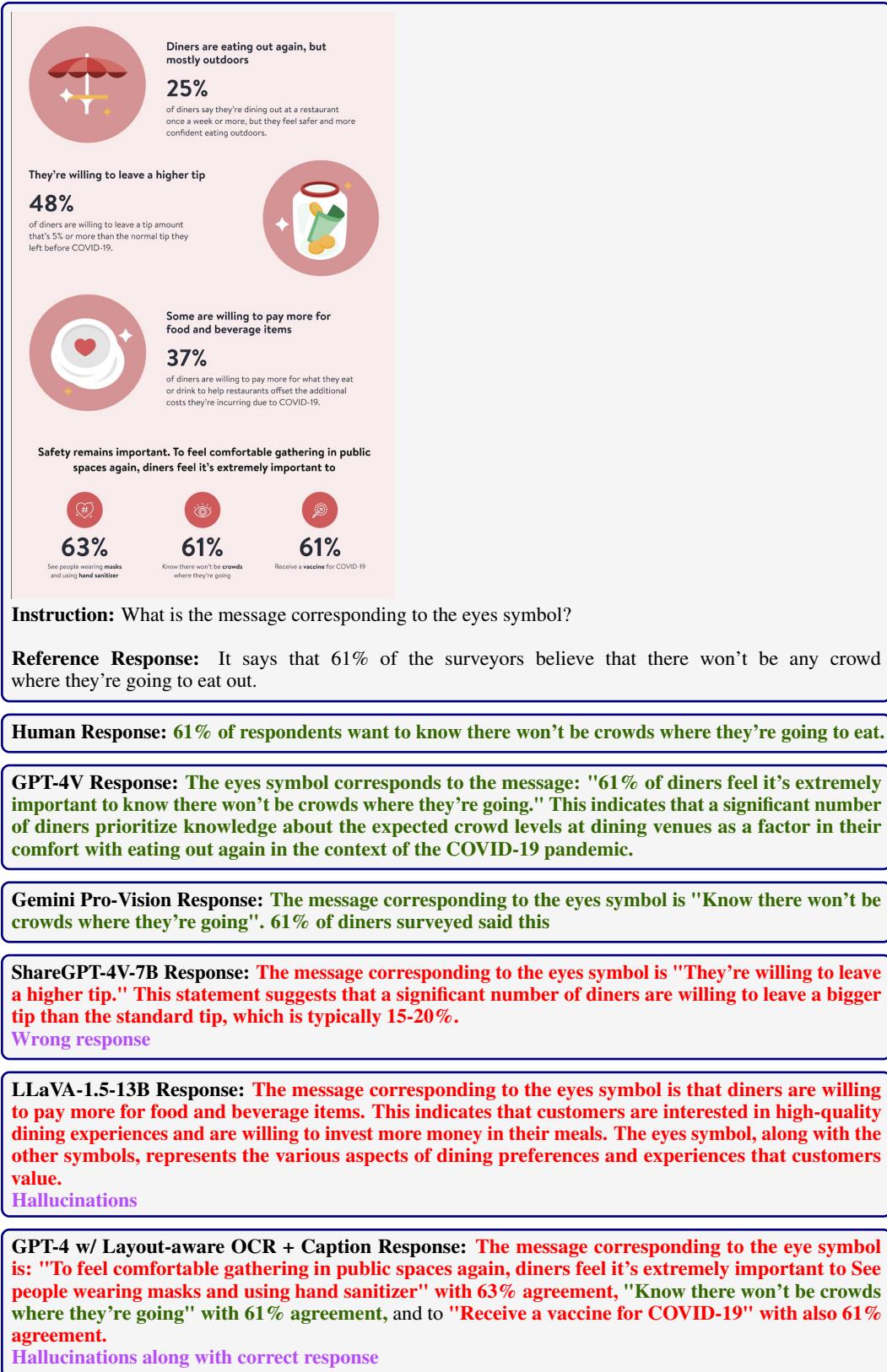


Figure 45: In this task, closed-source LMMs *GPT-4V* and *Gemini Pro-Vision* correctly respond to the instruction. This can be attributed to their ability for **visual perception, text recognition, and context-sensitive text-vision understanding**. All other models **hallucinate** their responses.

Stamp of Approval

During the campaign, President Obama argued that his election would help restore the image the rest of the world has of the United States. In the six months since his election, his approval ratings at home have slipped, though they remain high. Around the rest of the world, opinion is mixed. A recent study by WorldPublicOpinion.org asked people in 21 countries whether they had confidence that Obama would "do the right thing" when it came to world affairs. Here is a look at their answers.

Percent of respondents who have a lot of confidence or some confidence in Obama's ability to do the right thing in world affairs.

Percentage of respondents who have not too much confidence or no confidence at all in Obama's ability to do the right thing in world affairs.



PERCENTAGES MAY NOT ADD TO 100
DUE TO SOME RESPONDENTS GIVING NO ANSWER.

A COLLABORATION BETWEEN GOOD AND META & MICHAEL NEWHOUSE

Instruction: Calculate the difference in the percent of respondents who have confidence and do not have confidence in Obama from the country with a doll in their stamp.

Reference Response: Country is Poland. Percent of respondents who approve is 47 and the percent who do not have confidence is 38. The difference of the two is 9 percent.

Human Response: The country is Poland and the difference between the respondents who have and do not have confidence in Obama is 9%.

GPT-4V Response: To find the country with a doll in its stamp, I'll first identify that country and then calculate the difference in percentage of respondents who have confidence and do not have confidence in Obama from that country. Upon inspection, the country represented by a stamp with a doll on it is Poland. According to the data provided: 47% of respondents have a lot of confidence or some confidence in Obama's ability to do the right thing in world affairs. 38% of respondents do not have too much confidence or no confidence at all in Obama's ability to do the right thing in world affairs. To calculate the difference, we subtract the smaller percentage from the larger percentage: 47% (Confidence) - 38% (No confidence) = 9% Therefore, the difference in the percent of respondents who have confidence and do not have confidence in Obama from Poland is 9%.

Gemini Pro-Vision Response: For the country with the doll stamp (Nigeria), 85% of respondents have confidence in Obama, and 13% do not have confidence in Hallucinations

ShareGPT-4V-7B Response: The country with the highest percentage of people with confidence in Obama is Great Britain, with 80% of the respondents having confidence. On the other hand, the country with the lowest percentage of people with confidence in Obama is Pakistan, with only 20% of the respondents having confidence. The difference in the percentages between two countries is 60%. Hallucinations

LLaVA-1.5-13B Response: In the image, the country with a doll in their stamp is the United States. The percentage of respondents who have confidence in Obama is 63%, while the percentage of respondents who do not have confidence in Obama is 37%. The difference between two percentages is $63 - 37 = 26\%$. Hallucinations

GPT-4 w/ Layout-aware OCR + Caption Response: The image provided, unfortunately, doesn't contain information showing a specific country with a doll in their stamp. Hallucinations

Figure 46: In this task, **GPT-4V correctly** responds to the instruction. This can be attributed to their ability for **math and context-sensitive text-vision understanding**. All other models **hallucinate** their responses.

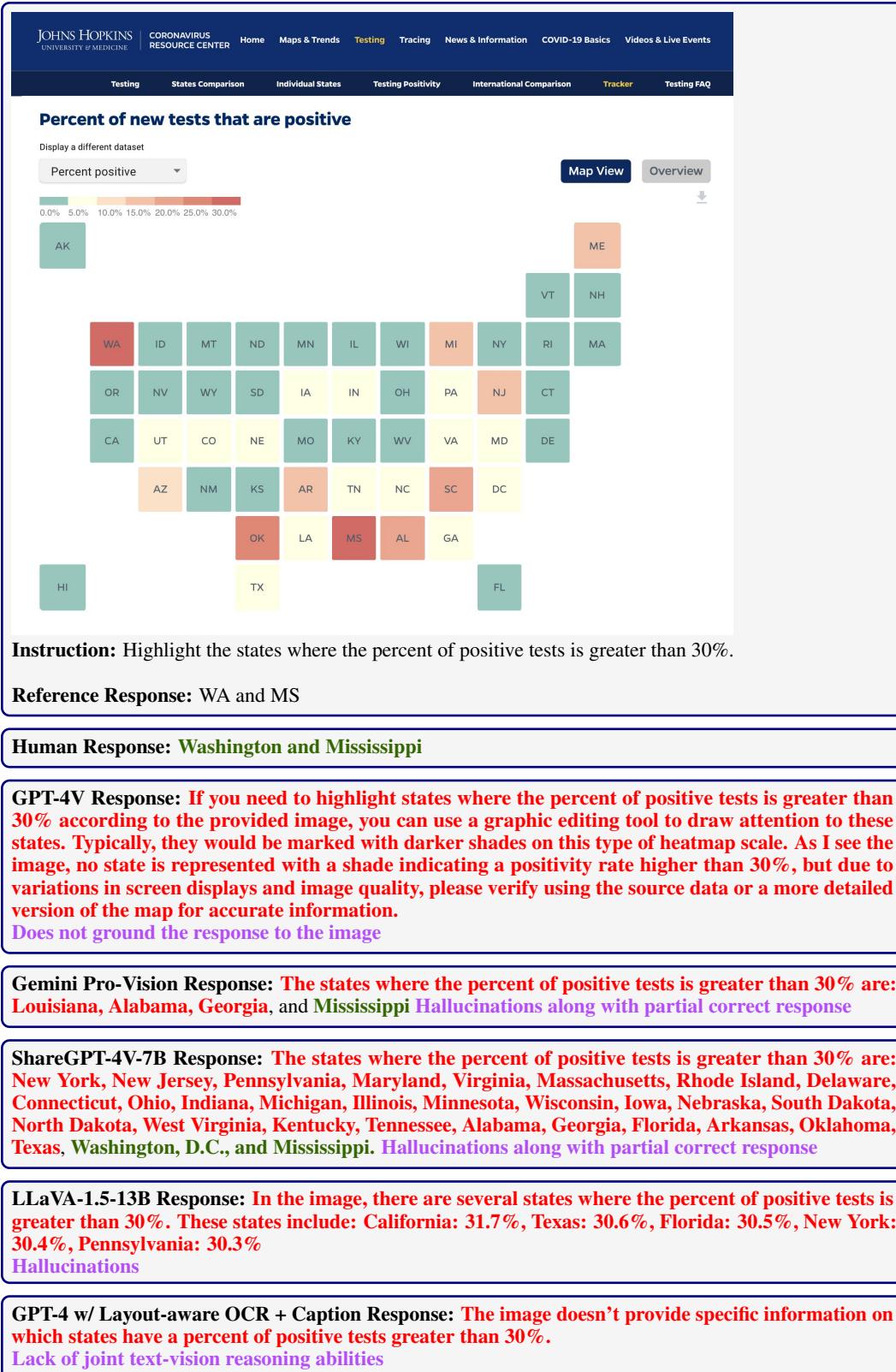


Figure 47: In this task, the closed source LMMs, **GPT4V does not ground the response** to the image. **Gemini-Vision-Pro, ShareGPT-4V-7B & LLaVA-1.5-13B** hallucinate in their responses. **GPT-4 w/ Layout-aware OCR + Caption Response** (Augmented LLM approach) **produces the wrong answer**, possibly due to lack of joint text-vision reasoning.



Instruction: How many production methods are used under pour-over/drip category?

Reference Response: 10

Human Response: There are 10 methods in the pour-over/drip category.

GPT-4V Response: In the pour-over/drip category of this coffee chart, there appear to be six production methods used. These methods are as follows: Auto Drip, Cone, Flat Bottom, AeroPress, Woodneck, Chemex

Wrong response

Gemini Pro-Vision Response: There are 8 production methods listed under pour-over/drip category.
Wrong response

ShareGPT-4V-7B Response: There are two production methods used under the pour-over/drip category: French press and drip coffee maker.
Wrong response

LLaVA-1.5-13B Response: There are three production methods used under the pour-over/drip category.

Wrong response

GPT-4 w/ Layout-aware OCR + Caption Response: The OCR text does not provide specific information on the number of production methods used under the pour-over/drip category.
Lack of joint text-vision reasoning

Figure 48: In this task, the *all LMMs*, produce the **wrong** response due to lack of fine-grained visual perception and possibly poor math reasoning. *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) **produces the wrong answer**, possibly due to lack of joint text-vision reasoning.



Instruction: Provide your opinion on whether changing the switch's current position will activate the instrument.

Reference Response: Despite turning the switch to the On position, the system will remain powered off due to the absence of a connected power cable.

Human Response: It will not. The instrument needs to be plugged in to work.

GPT-4V Response: The image shows the back panel of an electronic device with a power switch in the 'OFF' position. Changing the switch's current position from 'OFF' to 'ON' is likely to activate the instrument, provided that it is properly connected to a power source and there are no malfunctions with the instrument itself.

Wrong response

Gemini Pro-Vision Response: Flipping the switch to the "ON" position will power the instrument.
Wrong response

ShareGPT-4V-7B Response: Based on the image, it is not possible to determine if changing the switch's current position will activate the instrument. The switch is currently in the off position, and the instrument is not in use. To determine if changing the switch's position will activate the instrument, one would need to observe the instrument's response to the switch's movement.

Wrong response

LLaVA-1.5-13B Response: Based on the image, the switch is currently in the "off" position, which means that the instrument is not activated. To activate the instrument, the switch would need to be moved to the "on" position. This would allow the electrical current to flow through the device, enabling it to function properly.

Wrong response

GPT-4 w/ Layout-aware OCR + Caption Response: The instruction is ambiguous with not enough visual information available in the provided OCR text or the image caption regarding a switch. Therefore, it is not possible to provide an opinion on whether changing the switch's current position will activate the instrument.

Lack of joint text-vision reasoning

Figure 49: In this task, *GPT4V* produces the **correct** response, whereas *all other LMMs*, produce the **wrong** response due to lack of fine-grained visual perception. *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) **produces the wrong answer**, possibly due to lack of joint text-vision reasoning.



Instruction: What is written on a leg of the letter M?

Reference Response: SOHO

Human Response: SOHO

GPT-4V Response: On the leg of the letter "M" in the image, it says "SOHO."

Gemini Pro-Vision Response: SOHO

ShareGPT-4V-7B Response: A star is written on the leg of the letter M.

Wrong response

LLaVA-1.5-13B Response: The word "Music" is written on the leg of the letter M.

Wrong response

GPT-4 w/ Layout-aware OCR + Caption Response: U

Wrong response

Figure 50: In this task, **GPT-4V** and **Gemini Pro-Vision** **correctly** respond to the instruction. This can be attributed to their ability for text recognition, and context-sensitive text-vision understanding. Open source LMMs like **ShareGPT4V** and **LLaVA-V-1.5** and **GPT-4 w/ Layout-aware OCR + Caption Response** (Augmented LLM approach) predict the **wrong** answer.



Instruction: What would the first line read if we filter the letters in the red?

Reference Response: SAYRCON

Human Response: SA YR CON

GPT-4V Response: The first line reads "SATYRICONPDX.COM" with the letters in red. If we filter out the letters in red, the remaining letters spell "SATICONDX.COM".

Wrong response

Gemini Pro-Vision Response: SATYRICONPDX.COM

Wrong response

ShareGPT-4V-7B Response: If we filter the letters in the red, the first line would read "SATYRICONP"

Wrong response

LLaVA-1.5-13B Response: If we filter the letters in the red, the first line would read "Satyricon PDX".

Wrong response

GPT-4 w/ Layout-aware OCR + Caption Response: The instruction can't be executed on the provided information because the layout aware OCR text doesn't indicate any letters in red.

Lack of joint text-vision reasoning

Figure 51: In this task, *all LMMs* produce the wrong answer due to lack of fine-grained visual perception and context sensitive text-vision reasoning. *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) **produces the wrong answer**, possibly due to lack of joint text-vision reasoning.



Instruction: Get the number of the boat with three yellow and one red round items hanging from it.

Reference Response: WH97

Human Response: 97

GPT-4V Response: The boat with three yellow and one red round items (which are likely buoys) hanging from it has the number SS 170.
Hallucinations

Gemini Pro-Vision Response: SS 273

Wrong response

ShareGPT-4V-7B Response: The boat with three yellow and one red round items hanging from it is SS138.
Wrong response

LLaVA-1.5-13B Response: The boat with three yellow and one red round items hanging from it is the one with the number "S527" on it.
Hallucinations

GPT-4 w/ Layout-aware OCR + Caption Response: The instruction doesn't specify which boat has three yellow and one red round items hanging from it, so it's impossible to provide the number of that boat based on the provided OCR text and image caption.
Lack of joint text-vision reasoning

Figure 52: In this task, *GPT4V* & *LLaVA-1.5-13B* hallucinate, whereas, *Gemini-Vision-Pro* & *ShareGPT-4V-7B* produce the wrong response. The mispredictions of LMMs is possibly due to lack of fine-grained visual perception. *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) produces the wrong answer, possibly due to lack of joint text-vision reasoning.