

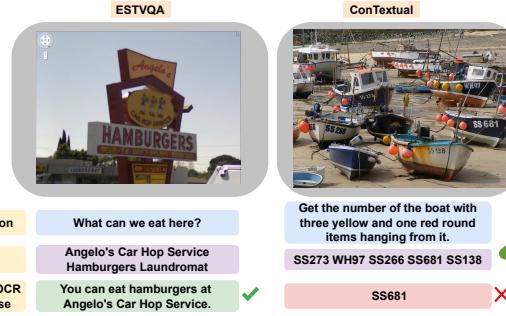
CONTEXTUAL: Evaluating Context-Sensitive Text-Rich Visual Reasoning in Large Multimodal Models

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

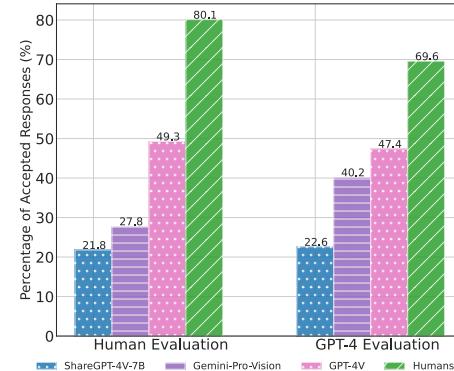
IMP: added images that cannot be solved simply by OCR as a meta-data to the prompt

Abstract

Many real-world tasks require an agent to reason jointly over text and visual objects, (e.g., navigating in public spaces), which we refer to as context-sensitive text-rich visual reasoning. Specifically, these tasks require an understanding of the context in which the text interacts with visual elements within an image. However, there is a lack of existing datasets to benchmark the state-of-the-art multimodal models' capability on context-sensitive text-rich visual reasoning. In this paper, we introduce CONTEXTUAL, a novel dataset featuring human-crafted instructions that require context-sensitive reasoning for text-rich images. We conduct experiments to assess the performance of 14 foundation models (GPT-4V, Gemini-Pro-Vision, LLaVA-Next) and establish a human performance baseline. Further, we perform human evaluations of the model responses and observe a significant performance gap of 30.8% between GPT-4V (the current best-performing Large Multimodal Model) and human performance. Our fine-grained analysis reveals that GPT-4V encounters difficulties interpreting time-related data and infographics. However, it demonstrates proficiency in comprehending abstract visual contexts such as memes and quotes. Finally, our qualitative analysis uncovers various factors contributing to poor performance including lack of precise visual perception and hallucinations. Our dataset, code, and leaderboard can be found on the project page <https://con-textual.github.io/>.



(a) Comparing an instance of CONTEXTUAL to existing datasets (e.g., ESTVQA). CONTEXTUAL requires contextualized understanding of the interactions between the textual and visual elements in the image while ESTVQA can be solved solely through text-based reasoning combined with accurate OCR detection.



(b) Performance of GPT-4V, Gemini-Pro-Vision, ShareGPT-4V-7B, and humans on the CONTEXTUAL dataset, with Human Evaluation (left sub-graph) and GPT4 Evaluation (right sub-graph).

Figure 1. Comparisons between our dataset CONTEXTUAL and prior work ESTVQA, along with benchmark performances of large multimodal models on CONTEXTUAL.

1. Introduction

The recent development of large multimodal models (LMMs) has resulted in models capable of responding to human instructions, posed as questions or imperative tasks, over images (Liu et al., 2023b; Chen et al., 2023; OpenAI, 2023b; Team et al., 2023; Dai et al., 2023; Bai et al., 2023;

^{*}Equal contribution ¹Department of Computer Science, University of California Los Angeles, USA. Correspondence to: Rohan Wadhawan <rwdhawan7@g.ucla.edu>, Hritik Bansal <hbansal@g.ucla.edu>.

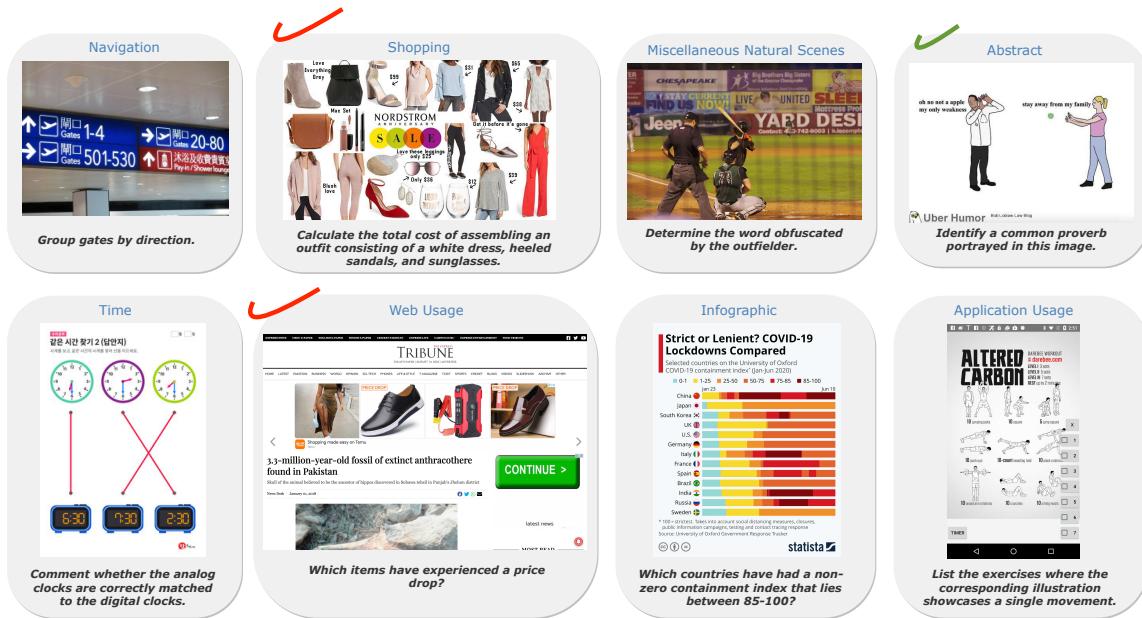


Figure 2. A sample (image, instruction) from each of the 8 visual scenarios in CONTEXTUAL dataset. The categories organized in a left-to-right, top-to-bottom reading order include Navigation, Shopping, Miscellaneous Natural Scenes, Abstract, Time, Web Usage, Infographic, and Application Usage.

HuggingFace, 2023; Ye et al., 2023a). Many real-world images contain texts within them which provides cues for comprehensively understanding them. 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).

Previous datasets such as TextVQA (Singh et al., 2019), STVQA (Singh et al., 2019), ESTVQA (Wang et al., 2020) have been proposed to assess the visual reasoning ability of multi-modal models over text-rich images. However, these datasets focused on accessing the OCR capability of the models to *read* the text in the image, and they usually do not require the model to capture the visual context in the image to answer the question. For example, in Figure 1a, we highlight an example from the ESTVQA dataset. Here, we show that a high accuracy OCR of the images (e.g., ‘Angelo’s Car Hop Service Hamburgers Laundromat’) has sufficient signal to answer the question (e.g., ‘What can we eat here?’). Though accessing the OCR capability is important, these examples do not test the unique potential of the LMMs to jointly reason over the embedded text and visual context in the image.

To evaluate multimodal models’ capability of jointly rea-

soning over embedded text and visual context in text-rich images, we propose CONTEXTUAL, a Context-sensitive Text-rich visual reasoning dataset consisting of 506 challenging instructions for LMMs evaluation. CONTEXTUAL covers eight real-world scenarios with text-rich images: time reading, shopping, navigation, abstract scenes, mobile application, webpages, infographics, and miscellaneous natural scenes (Figure 2). The diverse visual nature of these categories enables us to conduct a detailed, nuanced evaluation of the model’s capabilities.

Each instance of CONTEXTUAL contains a human-written instruction (question or imperative task) and a human-written ground-truth response (§2), with the constraint that to respond to an instruction accurately, a model must require context-sensitive joint reasoning over the textual and visual cues in the image. Figure 1a shows an example from our dataset. The instruction (‘Get the number of the boat with three yellow and one red round items hanging from it.’) cannot be answered even by perfectly capturing the OCR of the text content within the image (e.g., ‘SS273 WH97 SS266 SS681 SS138’). We summarize our work compared to the related works in Table 1.

We conduct extensive experiments using CONTEXTUAL to assess the reasoning abilities of 14 foundation models over context-sensitive text-rich images (§3.1). This includes three augmented LLMs setups (e.g., GPT-4 (OpenAI, 2023a)

IMP: GAP: VQA was done on text rich dataset which utilises the OCR as meta-data

Table 1. Comparison with related works for evaluating large multimodal models. We abbreviate Context-sensitive as Consens., Generation as Gen. We compare our work with LLaVA (Liu et al., 2023b), VisIT (Bitton et al., 2023), (Singh et al., 2019), STVQA (Biten et al., 2019), DUDE (Van Landeghem et al., 2023), InfoVQA (Mathew et al., 2022), and SEED Bench (Li et al., 2023a).

	Ours	LLaVA	VisIT	TextVQA	STVQA	DUDE	InfoVQA	SEED
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	✓	✗	✓	✗	✗	✓	✗	✓

prompted with combinations of image OCR, image layouts, and image captions), two proprietary LMMs (e.g., GPT-4V(OpenAI, 2023b), Gemini-Pro-Vision (Team et al., 2023)), and nine open LMMs (e.g., LLaVA-Next (Liu et al., 2024), LLaVA-1.5 (Liu et al., 2023a), ShareGPT-4V(Chen et al., 2023), Idefics (HuggingFace, 2023)). 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, §3.3).

Through human evaluations 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 1b). 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, only 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 (Sellam et al., 2020)) on the complete dataset for all the models (§3.3.1). Further, we perform fine-grained experiments to assess the model’s performance across visual contexts (§3.4). We observe that GPT-4V, the best performing LMM, struggles with time reading and infographic visual contexts, except for abstract contexts like

memes and quotes, where it outperforms humans. On the other hand, open models lag behind proprietary ones across most visual tasks, showing moderate proficiency only in abstract and natural scenes, owing to the need for more diversity of visual context in training data. However, we observe significant improvement in model performance with enhancement in image encoders, as seen with LLaVA-Next over LLaVA-v1.5. Lastly, we conduct a qualitative analysis (§4) of model responses for the different visual contexts in CONTEXTUAL, revealing that both proprietary and open models exhibit a limited capacity for fine-grained visual perception, with open models performing worse.

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 1 (b)). Our dataset encompasses a variety of tasks across diverse natural and digital text-rich visual scenarios, enabling robust testing.

Broadly, our benchmark follows these key collection guidelines: (a) Each sample consists of an $\langle \text{image}, \text{instruction}, \text{response} \rangle$ 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. (b) We aim to cover diverse instructions, including questions and imperative tasks. This ensures that the resulting dataset demonstrates a rich variety of instructions. (c) We aim to create instructions of varied complexity.

how to
perform
automatic
evaluation
??

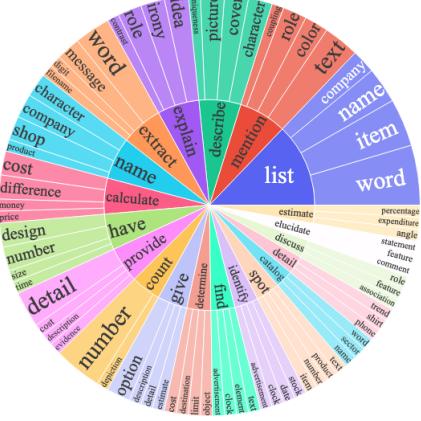


Figure 3. Top 20 Most frequently occurring verbs (inner circle) and their top 4 direct nouns (outer circle) in the instructions.

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.’).

Statistic	Number
Total number of samples	506
Root verbs in instructions	79
# Visual Contexts	8
- 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.

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, navigation, time, mobile and web usage, infographics, abstract scenes, and miscellaneous natural scenes. More details are available in Appendix §A.

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 (Schuhmann et al., 2022). Specifically, we use keyword search using CLIP-retrieval UI (Beaumont, 2022). 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., gro-

cery, 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, calendars), and Abstract (e.g., memes, quotes, comic strips, science jokes, math jokes).

Secondly, we source images for the *Application Usage* category from the Rico Dataset (Deka et al., 2017), which includes 66,000 distinct UI screens originating from 9,300 Android apps. Thirdly, we scrape the website links made available by the Open WebText Initiative (Gokaslan & Cohen, 2019) and collect screenshots for the *Web Usage* category. Lastly, we acquire images from the test sets of existing VQA datasets, and proceed to annotate them with novel instruction-response pairs. Originally, these datasets consists question-and-answer pairs that primarily focus on text recognition capabilities. However, they offer an opportunity to formulate context-sensitive instructions for the images. Specifically, we reannotate these image instances, categorizing them into two groups: *Infographic*, sourced from the InfographicVQA (Mathew et al., 2022) dataset, and *Miscellaneous Natural Scenes*, sourced from the STVQA (Biten et al., 2019) and ESTVQA (Wang et al., 2020) datasets.

2.3. Data Annotation

Stage 1: 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 (padlepadle, 2023) 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: We divide 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.¹

Stage 3: In this final stage, we perform a verification process for each sample annotated in Stage 2. We asked MTurk workers (mutually exclusive from the ones used for human performance baseline and human evaluation) to verify the correctness of each sample *<image, instruction, response>* and found that 96% of the samples were annotated correctly. Filtering out the incorrect samples, we tasked each author group to review the samples created by the other group. This ensured adherence to guidelines and filtered out low-quality samples. Finally,

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

we end up with a dataset of **506** instances.

We provide statistics for 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 3. We also annotate each sample to determine whether it requires information extraction, and mathematical reasoning (Appendix §G.1). We provide details on data release in §C.

3. Experiments

3.1. Setup

Augmented LLMs. 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 (Lu et al., 2023b; Wu et al., 2023; Surís et al., 2023; Gupta & Kembhavi, 2023). 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 (paddlepadle, 2023) for extracting OCR from the images, **LATIN** prompt (Wang et al., 2023a) inspired OCR text arrangement implementation to maintain **layout-awareness in the OCR**, and **ShareGPT-4V-7B** for the dense image captions (App. §E).

LMMs. We evaluate GPT-4V (OpenAI, 2023b) and Gemini-Pro-Vision (Team et al., 2023) that are representative proprietary LMMs that have achieved state-of-the-art on other visual reasoning benchmarks (Goyal et al., 2017). In addition, we evaluate a wide range of open LMMs including **LLaVA-Next-34B** (Liu et al., 2024), **LLaVA-1.5-13B** (Liu et al., 2023a), **ShareGPT-4V-7B** (Chen et al., 2023), **mPLUG-Owl-v2-7B** (Ye et al., 2023a;b), **Qwen-VL-7B** (Bai et al., 2023), **InstructBLIP-Vicuna-7B** (Dai et al., 2023), and **Idefics-9B** (HuggingFace, 2023). We include LLaVAR (Zhang et al., 2023) and BLIVA-Vicuna-7B (Hu et al., 2023) as they were introduced for text-rich visual reasoning.

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 §B.1. We spent \$180 on collecting human predictions on our dataset.

3.2. Evaluation

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. We provide the screenshot of our annotation interface in Appendix B.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** (Sellam et al., 2020), Rouge-L (Lin, 2004) and BERTScore (Zhang et al., 2019) 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 §F. 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).

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

types of
OCR

Table 3. Comparison of the performance of various foundation models (augmented LLM and LMMs) and humans on a subset of CONTEXTUAL dataset (280 samples). 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. 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 (OpenAI, 2023b)	<u>49.3</u>	<u>47.4</u>	<u>45.0</u>	<u>45.3</u>	17.3	52.5
Gemini-Pro-Vision (Team et al., 2023)	27.8	40.2	37.1	42.5	<u>30.1</u>	<u>58.4</u>
LLaVA-1.5-13B (Liu et al., 2023a)	17.2	20.6	17.5	43.6	21.7	54.8
ShareGPT-4V-7B (Chen et al., 2023)	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

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

Table 4. Comparison of the human and automatic evaluation metric using ROC-AUC and spearman correlation on a subset of Contextual dataset (280 samples, similar to Table 3). We find that the GPT-4 and GPT-4V based evaluation correlate the most with humans.

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 outperform those from GPT-4V by 22.2% and 23.6% using GPT-4 and GPT-4V evaluation. Like 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 of GPT-4V and Gemini-Pro-Vision is 7.2% as per GPT4 evaluation**. We still observe a large gap in the 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.

Furthermore, we find that the BLEURT scores for humans

are the highest, while GPT-4V achieves the highest score among the LMMs. Interestingly, 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 ([Sellam et al., 2020](#)).

3.3.1. CORRELATION ANALYSIS

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 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. 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 ([Bitton et al., 2023](#)) that shows GPT-4 evaluation outperforms standard text generation metrics. Further, the dataset size is sufficient to get reliable confidence intervals on GPT-4 evaluation. We compared model predictions pairwise using the paired t-test at a 95% confidence interval. The comparison between LLaVA-v1.5 and GPT4v/Gemini-Pro-Vision yielded a P value < 0.0001, suggesting that the difference in the performance is statistically significant. Comparing GPT4V with Gemini-Pro-Vision resulted in a P value of

Table 5. Fine-grained comparison in the zero-shot performance of the foundation models and humans on the CONTEXTUAL dataset using GPT-4 evaluation. We abbreviate the average response acceptance rating as Avg., Navigation as Nav., Abstract as Abs., Application usage as App., 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 (OpenAI, 2023b)	47.4	<u>18.0</u>	<u>54.0</u>	<u>48.0</u>	100.0	48.0	42.0	<u>28.0</u>	48.0
Gemini-Pro-Vision (Team et al., 2023)	40.2	16.0	32.7	28.6	65.3	44.9	43.8	20.0	<u>52.8</u>
LLaVA-Next-34B (Liu et al., 2024)	36.8	10.0	36.0	30.6	66.0	36.0	28.0	12.0	51.3
ShareGPT-4V-7B (Chen et al., 2023)	22.6	0.0	16.0	20.0	28.6	20.0	20.0	14.0	37.7
Qwen-VL-7B (Bai et al., 2023)	21.8	4.0	20.0	24.0	53.1	6.0	18.0	14.0	27.3
LLaVA-1.5B-13B (Liu et al., 2023a)	20.8	4.0	10.0	18.0	44.9	16.0	26.0	4.0	29.7
mPLUG-Owl-v2-7B (Ye et al., 2023a)	18.6	4.0	8.0	24.0	32.7	20.0	10.0	12.0	26.0
LLaVAR-13B (Zhang et al., 2023)	14.9	10.0	16.0	6.0	44.9	8.0	10.0	6.0	16.7
BLIVA-Vicuna-7B (Hu et al., 2023)	10.3	2.0	4.0	14.0	24.5	4.0	8.0	4.0	14.7
InstructBLIP-Vicuna-7B (Dai et al., 2023)	9.7	2.0	4.0	16.0	20.0	6.0	12.0	2.1	12.0
Idefics-9B (HuggingFace, 2023)	7.7	4.0	2.0	12.0	12.0	0.0	6.0	2.0	13.3
Humans	69.6	64.0	64.0	73.5	75.5	64.0	58.0	72.0	78.0

0.035, also denoting statistical significance. Therefore, the differences in model performance on ConTextual are statistically significant at the 95% confidence level. As a result, we utilize GPT-4 for automatically evaluate the quality of the predicted responses.

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 ([Lu et al., 2023a; Masry et al., 2023](#)) has shown that the existing LMMs underperform humans in reasoning over charts.

Further, we observe that the best performing open model LLaVA-Next-34B bridges the gap between the other open source models like LLaVA-1.5-13B and ShareGPT-4V-7B and the closed source models like Gemini-Pro-Vision. It performs the best on abstract and natural scenes, while it struggles the most on time and infographics. The relative imbalance in performance across categories can be attributed

to the lack of diverse visual contexts in their training data. For instance, the COCO dataset ([Lin et al., 2014](#)) used for vision-language alignment in the open models predominantly comprises natural scenes. However, comparing it to its predecessor LLaVA-1.5-13B, improvement in visual encoding, data diversity, and LLM capacity boost performance on CONTEXTUAL. We also observe the open models specifically introduced for text-rich visual reasoning like LLaVAR and BLIVA-Vicuna-7B 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 §G. Overall, our fine-grained analysis aids in identifying the gaps in the existing models which would inspire the development of next generation LMMs.

3.5. Study on Synthetically Scaling Data

Creating synthetic data for context-sensitive text-rich visual reasoning is challenging. Automatic dataset generation using OCR and image caption data with LLMs like GPT-4, exemplified by LLaVAR ([Zhang et al., 2023](#)), yields instructions solvable by OCR+LLM or basic object understanding but shows poor performance on the context-sensitive instructions, as reported in Table 5. Further, finding suitable images for joint reasoning is complex, requiring heuristic OCR filtering and manual selection. Despite these challenges, a synthetic data generation pipeline has been developed, showing model performance on synthetic data similar to human-annotated data, demonstrating the potential to scale context-sensitive instructions effectively (refer to §H for more details).

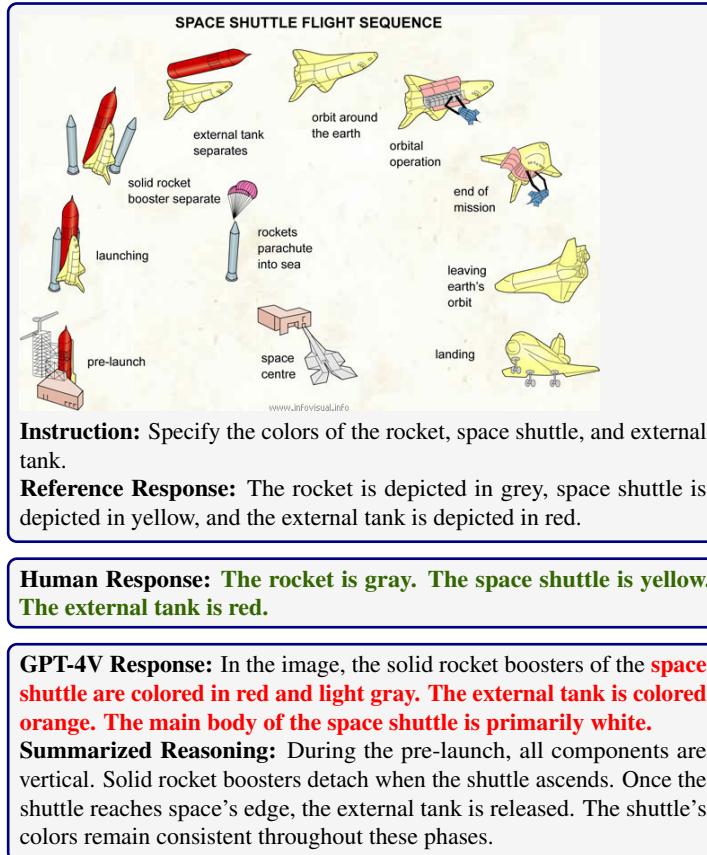


Figure 4. Here, GPT-4V responds incorrectly to the instruction despite its logical reasoning. Green denotes correct response, and red highlights response errors. *Summarized Reasoning:* rationale used by GPT-4V to answer.

4. Qualitative Examples

GPT-4V vs Humans. In Figure 4, we see an instance where GPT-4V provides an incorrect answer. Here, 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. We provide more examples in Appendix §I (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. A similar analysis was presented in the prior work (Guan et al., 2023) where GPT-4V fails on the perturbed versions of common visual illusions.

GPT-4V vs. Open LMMs and Augmented LLM. We compare the best-performing models in each category, closed-source LMM, open-source LMM, and Augmented LLM approach, that is, GPT-4V, LLaVA-Next-34B, and GPT-4 w/ Layout-aware OCR + Caption, respectively, using an example illustrated in Figure 5. GPT-4V correctly identi-

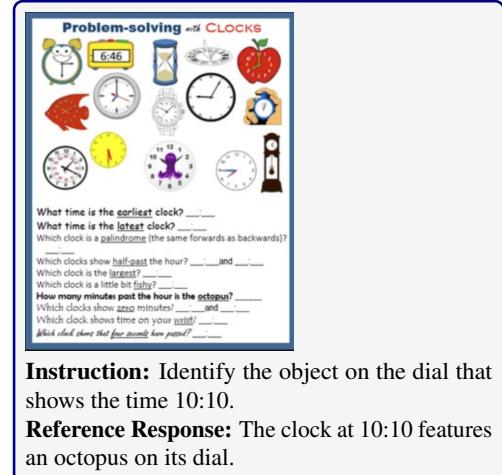


Figure 5. GPT-4V correctly responds to the instruction. However, LLaVA-Next-34B and GPT-4 w/ Layout-aware OCR+ Caption (Augmented LLM) produce wrong responses.

fies the object, showcasing superior visual perception and context-sensitive text-rich visual reasoning abilities over the LLaVA-Next and Augmented LLM approach that produces the wrong response. LLaVA-Next does not ground its response to the image due to relatively poor context-sensitive text-rich visual reasoning abilities. On the other hand, the Augmented LLM approach cannot respond to this instruction because the image caption and layout-aware OCR do not provide sufficient information to reason over embedded text and visual elements in the image. We refer to Figure 28, 31, 32, 38, 42, 45, 46, 50 for more examples demonstrating instances where open models exhibit lack of context-sensitive text-rich visual reasoning, or deficiencies in fine-grained perception.

Our analysis suggests that enhancing image encoders and increasing training data diversity can improve model perception, leading to more effective context-sensitive reasoning in text-rich visual contexts.

5. Related Work

Text-Rich Image Understanding. Recently, there has been a growing interest in understanding the interactions between the text and visual elements in the image (Lee et al., 2023; Xu et al., 2020). To track the progress of the models in this field, several datasets were introduced like OCRVQA (Mishra et al., 2019), TextVQA (Singh et al., 2019), DocVQA (Mathew et al., 2021), STVQA (Biten et al., 2019), ESTVQA (Wang et al., 2020). These datasets majorly focus on the ability of the models to accurately read the text in the documents or natural scene images. Prior work (Liu et al., 2023c) provides a benchmark to assess the ability of LMMs to perform accurate OCR. 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 (Antol et al., 2015; Goyal et al., 2017), image captioning (Gurari et al., 2018; Lin et al., 2014) where the model primarily needs to understand the key objects and their relations. Later, several works were introduced to assess the commonsense reasoning, which requires the models to reason about the questions that require skills beyond recognition, including VCR (Zellers et al., 2019; Yin et al., 2021). In addition, there are several datasets and benchmarks that evaluate specific skills of the LMMs including math skills (Chen et al., 2022; Lu et al., 2021a;b; 2023a), world knowledge (Yue et al., 2023), and grade school science diagrams (Kembhavi et al., 2016; Wang et al., 2023b). Additionally, there are several datasets for meme understanding such as hateful memes (Kiela et al., 2020), memefiy (Vyalla & Udandarao, 2020), and memecap (Hwang & Shwartz, 2023). Such works will require joint reasoning over text and visual content over the image. However, in our work, we broaden the scope and identify a breadth of visual domains that require context-sensitive text-rich visual reasoning. These include time reading, navigation and transportation in public spaces, meme and quote understanding, and shopping etc.

Large Multimodal Models. Prior works such as LXMERT (Tan & Bansal, 2019), VisualBERT (Li et al., 2019), X-decoder (Zou et al., 2023) learn robust vision-language representations by training on image-text data such as Conceptual captions (Changpinyo et al., 2021), COCO (Lin et al., 2014). Post-training, they will be finetuned on the specific tasks such as VQA (Antol et al., 2015). For doc-

ument understanding, popular vision-language models include Pix2Struct (Lee et al., 2023), Donut (Kim et al., 2022), MatCha (Liu et al., 2022). However, the development of large language models (Brown et al., 2020; OpenAI, 2023a), trained on internet-scale text corpus, shifted the paradigm towards the development of general-purpose multimodal models. Specifically, these are vision-language generative models that can solve diverse tasks in a zero-shot manner without task-specific finetuning. Notably, these are popularly known as large multimodal models (LMMs). These include proprietary models such as GPT-4V (OpenAI, 2023b) and Gemini-Pro-Vision (Team et al., 2023). These models have achieved state-of-the-art performance on the traditional vision-language models. In the open space, the models include LLaVA (Liu et al., 2023b;a; 2024), mPLUG-Owl (Ye et al., 2023a), OpenFlamingo (Awadalla et al., 2023), Idefics (HuggingFace, 2023), LLaMA-Adapter (Gao et al., 2023), Idefics (HuggingFace, 2023). In addition, there are a class of LMMs that focus on enhanced text-rich visual reasoning capabilities including LLAVAR (Zhang et al., 2023) and BLIVA (Hu et al., 2023). 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. We ask humans to solve our dataset and also use human annotators for model response evaluation. We find that the modern LMMs (proprietary and open models) struggle to perform on our dataset while humans are good at it. In summary, our dataset paves a path for assessing the progress on reasoning over text-rich images, a domain with significant real-world applications. We make the dataset² and code³ available to the LMM community along with a continuously updated leaderboard⁴ with recent LMMs.

Acknowledgements

This research is supported in part by the ECOLE program under Cooperative Agreement HR00112390060 with the

²Hugging Face Dataset

³GitHub Code Repository

⁴HuggingFace Leaderboard

US Defense Advanced Research Projects Agency (DARPA), and UCLA-Amazon Science Hub for Humanity and Artificial Intelligence. Hritik Bansal is supported in part by AFOSR MURI grant FA9550-22-1-0380.

Impact Statement

CONTEXTUAL is proposed to evaluate the context-sensitive text-rich visual reasoning capabilities of large multimodal models. These models are a class of generative models provide textual response to user instructions, grounded in text, for diverse images. During our data collection, we aim to ensure that the images, human-written instructions, and reference responses are not offensive to any social group. We are aware that the existing multimodal models are capable of generating harmful responses, despite the presence of safeguard filter. In addition, our qualitative analysis reveals that the model responses would hallucinate, however, we did not observe any apparent harmful and privacy sensitive information in them.

In our experiments, we asked human annotators, mainly from the US, to provide responses to establish a human baseline. We are aware that the linguistic diversity and writing style of the human responses would change with different social groups. The extension of our work can focus on understanding the impact of different social groups on the human baseline performance on the CONTEXTUAL dataset. A similar argument is relevant for human evaluation of the model responses. To obtain more reliable human evaluation results, future work would involve annotators from more diverse regions.

References

- Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C. L., and Parikh, D. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pp. 2425–2433, 2015.
- Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa, S., et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- Beaumont, R. Clip retrieval: Easily compute clip embeddings and build a clip retrieval system with them. <https://github.com/rom1504/clip-retrieval>, 2022.
- Biten, A. F., Tito, R., Mafla, A., Gomez, L., Rusinol, M., Valveny, E., Jawahar, C., and Karatzas, D. Scene text visual question answering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4291–4301, 2019.
- Bitton, Y., Bansal, H., Hessel, J., Shao, R., Zhu, W., Awadalla, A., Gardner, J., Taori, R., and Schimdt, L. Visit-bench: A benchmark for vision-language instruction following inspired by real-world use. *arXiv preprint arXiv:2308.06595*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Changpinyo, S., Sharma, P., Ding, N., and Soricut, R. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3558–3568, 2021.
- Chen, J., Li, T., Qin, J., Lu, P., Lin, L., Chen, C., and Liang, X. Unigeo: Unifying geometry logical reasoning via reformulating mathematical expression. *arXiv preprint arXiv:2212.02746*, 2022.
- Chen, L., Li, J., Dong, X., Zhang, P., He, C., Wang, J., Zhao, F., and Lin, D. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- Dai, W., Li, J., Li, D., Tiong, A. M. H., Zhao, J., Wang, W., Li, B., Fung, P., and Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- Deka, B., Huang, Z., Franzen, C., Hirschman, J., Afergan, D., Li, Y., Nichols, J., and Kumar, R. 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*, pp. 845–854, 2017.
- Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., Li, H., and Qiao, Y. Llama-adapter v2: Parameter-efficient visual instruction model, 2023.
- Gokaslan, A. and Cohen, V. Openwebtext corpus. <http://skylion007.github.io/OpenWebTextCorpus>, 2019.
- Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., and Parikh, D. 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*, pp. 6904–6913, 2017.

- Guan, T., Liu, F., Li, X. W. R. X. Z., Wang, X. L. X., Yacoob, L. C. F. H. Y., and Zhou, D. M. T. Hallusionbench: An advanced diagnostic suite for entangled language hallucination & visual illusion in large vision-language models. *arXiv e-prints*, pp. arXiv–2310, 2023.
- Gupta, T. and Kembhavi, A. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14953–14962, 2023.
- Gurari, D., Li, Q., Stangl, A. J., Guo, A., Lin, C., Grauman, K., Luo, J., and Bigham, J. P. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3608–3617, 2018.
- Hu, W., Xu, Y., Li, Y., Li, W., Chen, Z., and Tu, Z. Bliva: A simple multimodal llm for better handling of text-rich visual questions. *arXiv preprint arXiv:2308.09936*, 2023.
- HuggingFace. Introducing idefics: An open reproduction of state-of-the-art visual language model. <https://huggingface.co/blog/idefics>, 2023.
- Hwang, E. and Shwartz, V. Memecap: A dataset for captioning and interpreting memes. *arXiv preprint arXiv:2305.13703*, 2023.
- Kembhavi, A., Salvato, M., Kolve, E., Seo, M., Hajishirzi, H., and Farhadi, A. 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*, pp. 235–251. Springer, 2016.
- Kiela, D., Firooz, H., Mohan, A., Goswami, V., Singh, A., Ringshia, P., and Testuggine, D. The hateful memes challenge: Detecting hate speech in multimodal memes. *Advances in neural information processing systems*, 33: 2611–2624, 2020.
- Kim, G., Hong, T., Yim, M., Nam, J., Park, J., Yim, J., Hwang, W., Yun, S., Han, D., and Park, S. Ocr-free document understanding transformer. In *European Conference on Computer Vision*, pp. 498–517. Springer, 2022.
- Lee, K., Joshi, M., Turc, I. R., Hu, H., Liu, F., Eisenschlos, J. M., Khandelwal, U., Shaw, P., Chang, M.-W., and Toutanova, K. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *International Conference on Machine Learning*, pp. 18893–18912. PMLR, 2023.
- Li, B., Wang, R., Wang, G., Ge, Y., Ge, Y., and Shan, Y. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023a.
- Li, L. H., Yatskar, M., Yin, D., Hsieh, C.-J., and Chang, K.-W. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- Li, Y., Wang, L., Hu, B., Chen, X., Zhong, W., Lyu, C., and Zhang, M. A comprehensive evaluation of gpt-4v on knowledge-intensive visual question answering. *arXiv preprint arXiv:2311.07536*, 2023b.
- Lin, C.-Y. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Liu, F., Piccinno, F., Krichene, S., Pang, C., Lee, K., Joshi, M., Altun, Y., Collier, N., and Eisenschlos, J. M. Matcha: Enhancing visual language pretraining with math reasoning and chart derendering. *arXiv preprint arXiv:2212.09662*, 2022.
- Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023a.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023b.
- Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., and Lee, Y. J. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.
- Liu, Y., Li, Z., Li, H., Yu, W., Huang, M., Peng, D., Liu, M., Chen, M., Li, C., Jin, L., et al. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023c.
- Lu, P., Gong, R., Jiang, S., Qiu, L., Huang, S., Liang, X., and Zhu, S.-C. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. *arXiv preprint arXiv:2105.04165*, 2021a.
- Lu, P., Qiu, L., Chen, J., Xia, T., Zhao, Y., Zhang, W., Yu, Z., Liang, X., and Zhu, S.-C. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. *arXiv preprint arXiv:2110.13214*, 2021b.
- Lu, P., Bansal, H., Xia, T., Liu, J., Li, C., Hajishirzi, H., Cheng, H., Chang, K.-W., Galley, M., and Gao, J. Mathvista: Evaluating mathematical reasoning of

- foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023a.
- Lu, P., Peng, B., Cheng, H., Galley, M., Chang, K.-W., Wu, Y. N., Zhu, S.-C., and Gao, J. Chameleon: Plug-and-play compositional reasoning with large language models. *arXiv preprint arXiv:2304.09842*, 2023b.
- Masry, A., Kavehzadeh, P., Do, X. L., Hoque, E., and Joty, S. Unichart: A universal vision-language pretrained model for chart comprehension and reasoning. *arXiv preprint arXiv:2305.14761*, 2023.
- Mathew, M., Karatzas, D., and Jawahar, C. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 2200–2209, 2021.
- Mathew, M., Bagal, V., Tito, R., Karatzas, D., Valveny, E., and Jawahar, C. Infographiccvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1697–1706, 2022.
- Mishra, A., Shekhar, S., Singh, A. K., and Chakraborty, A. Ocr-vqa: Visual question answering by reading text in images. In *2019 international conference on document analysis and recognition (ICDAR)*, pp. 947–952. IEEE, 2019.
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023a, 2023a.
- OpenAI. Gpt-4v(ision) system card, 2023b. <https://openai.com/research/gpt-4v-system-card>, 2023b.
- paddlepadle. Paddleocr: Multilingual ocr toolkit based on paddlepaddle. <https://github.com/PaddlePaddle/PaddleOCR>, 2023.
- Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., 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.
- Sellam, T., Das, D., and Parikh, A. P. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*, 2020.
- Singh, A., Natarajan, V., Shah, M., Jiang, Y., Chen, X., Batra, D., Parikh, D., and Rohrbach, M. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8317–8326, 2019.
- Surís, D., Menon, S., and Vondrick, C. Vipergrpt: Visual inference via python execution for reasoning. *arXiv preprint arXiv:2303.08128*, 2023.
- Tan, H. and Bansal, M. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., Soriceut, R., Schalkwyk, J., Dai, A. M., Hauth, A., et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Van Landeghem, J., Tito, R., Borchmann, Ł., Pietruszka, M., Joziak, P., Powalski, R., Jurkiewicz, D., Coustaty, M., Anckaert, B., Valveny, E., et al. Document understanding dataset and evaluation (dude). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 19528–19540, 2023.
- Vyalla, S. R. and Udundarao, V. Memefy: A large-scale meme generation system. In *Proceedings of the 7th ACM IKDD CoDS and 25th COMAD*, pp. 307–311. 2020.
- Wang, W., Li, Y., Ou, Y., and Zhang, Y. Layout and task aware instruction prompt for zero-shot document image question answering. *arXiv preprint arXiv:2306.00526*, 2023a.
- Wang, X., Liu, Y., Shen, C., Ng, C. C., Luo, C., Jin, L., Chan, C. S., Hengel, A. v. d., and Wang, L. 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*, pp. 10126–10135, 2020.
- Wang, X., Hu, Z., Lu, P., Zhu, Y., Zhang, J., Subramaniam, S., Loomba, A. R., Zhang, S., Sun, Y., and Wang, W. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*, 2023b.
- Wang, Y., Mishra, S., Alipoormolabashi, P., Kordi, Y., Mirzaei, A., Arunkumar, A., Ashok, A., Dhanasekaran, A. S., Naik, A., Stap, D., et al. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *arXiv preprint arXiv:2204.07705*, 2022.
- Wu, C., Yin, S., Qi, W., Wang, X., Tang, Z., and Duan, N. Visual chatgrpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- Xu, Y., Li, M., Cui, L., Huang, S., Wei, F., and Zhou, M. 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*, pp. 1192–1200, 2020.

- Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023a.
- Ye, Q., Xu, H., Ye, J., Yan, M., Liu, H., Qian, Q., Zhang, J., Huang, F., and Zhou, J. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.04257*, 2023b.
- Yin, D., Li, L. H., Hu, Z., Peng, N., and Chang, K.-W. Broaden the vision: Geo-diverse visual commonsense reasoning. *arXiv preprint arXiv:2109.06860*, 2021.
- Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.
- Zellers, R., Bisk, Y., Farhadi, A., and Choi, Y. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6720–6731, 2019.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., and Artzi, Y. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- Zhang, Y., Zhang, R., Gu, J., Zhou, Y., Lipka, N., Yang, D., and Sun, T. Llavar: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023.
- Zou, X., Dou, Z.-Y., Yang, J., Gan, Z., Li, L., Li, C., Dai, X., Behl, H., Wang, J., Yuan, L., et al. Generalized decoding for pixel, image, and language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15116–15127, 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.



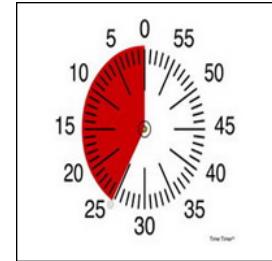
(a) Single Clock



(b) Multiple Clocks



(c) Calendar



(d) Timer

Figure 6. Examples of the Time visual Scenario



(a) Grocery



(b) Furniture



(c) Clothes



(d) Gadgets



(e) Cosmetics



(f) Services

Figure 7. Examples of the Shopping visual scenario

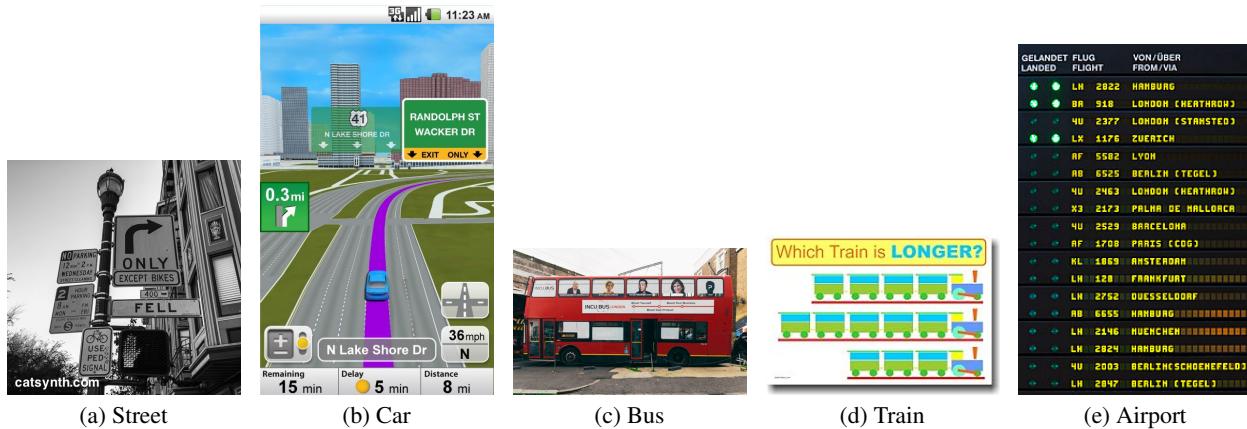
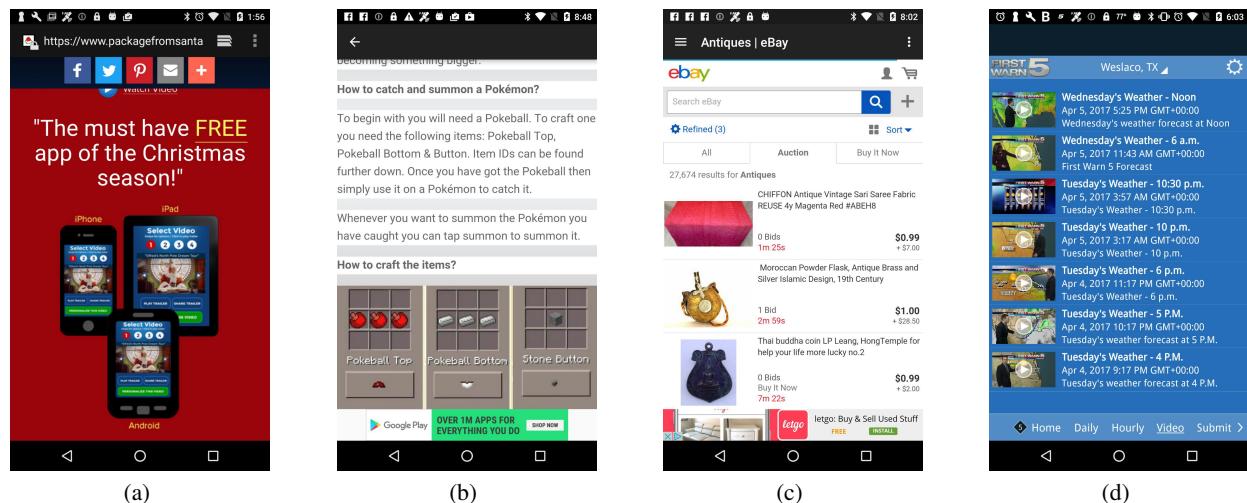

 Figure 8. Examples of the *Navigation* visual scenario

 Figure 9. Examples of the *Abstract* visual scenario

 Figure 10. Examples of the *Mobile Usage* visual scenario

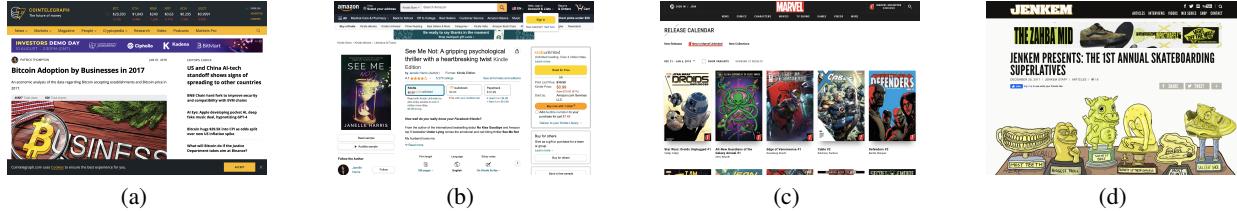


Figure 11. Examples of the *Web Usage* visual scenario

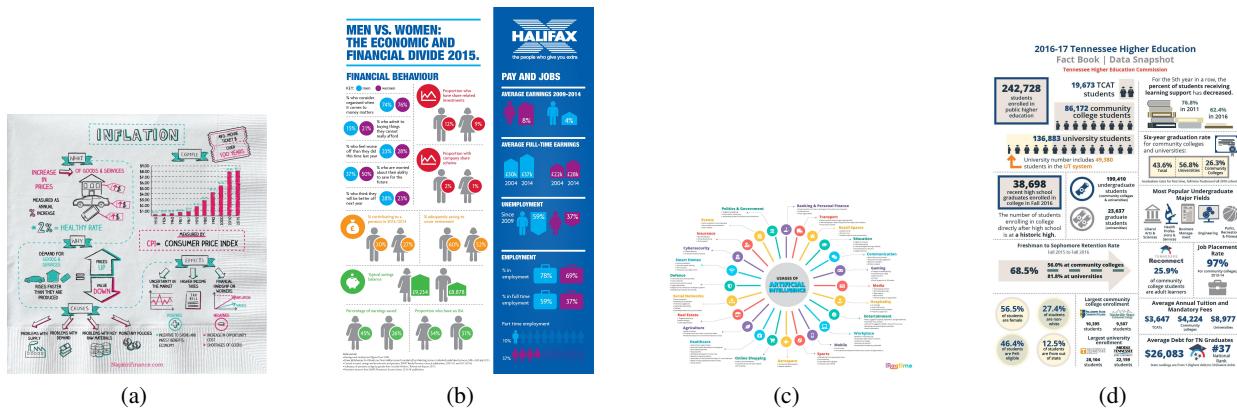


Figure 12. Examples of the *Infographic* visual scenario



Figure 13. Examples of the *Miscellaneous Natural Scenes* visual scenario

B. Human Annotation Screenshots

B.1. Human Performance Screenshot

We present the screenshot of the user interface used for acquiring human responses on the CONTEXTUAL dataset in Figure 14.

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 14. User interface of the human response collection.

B.2. Human Evaluation Screenshot

We present the screenshot of the user interface used for human evaluation in Figure 15.

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 15. User interface of the human evaluation.

C. 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.

D. 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 16. 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.

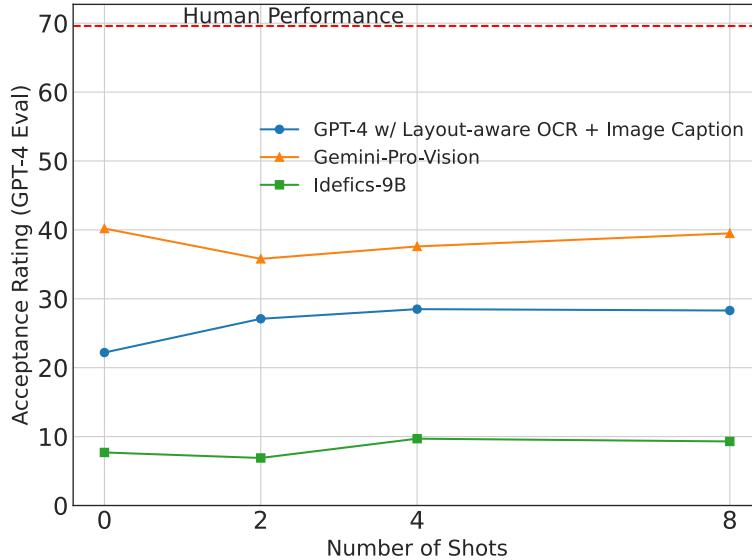


Figure 16. Few-shot performance on the CONTEXTUAL dataset.

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 (Li et al., 2023b). 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.

E. 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.

E.1. GPT-4 w/ OCR

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 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. 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 Response: Both the Living Room and the Hip Living Room include 4 items.

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

E.2. GPT-4 w/ Layout-aware OCR



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 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.

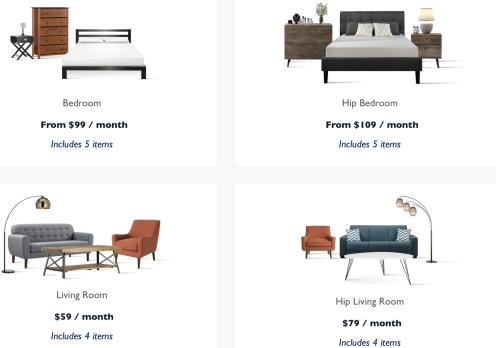
OCR TEXT:

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 GPT-4 w/ Layout-aware OCR Response

E.3. GPT-4 w/ Layout-aware OCR + Caption Response



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 GPT-4 w/ Layout-aware OCR + Caption Response.

F. GPT-4 Evaluation Prompt

In this section, we provide an illustration of a GPT4 prompt used to assess both model-generated and human responses. Figure 20 displays an instance within the Shopping category, featuring a reference response generated by the Gemini-Pro-Vision model. It's important to observe that the prompt does not include any information about the predicting model.



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 ImageTaskEvaluatorGPT, an expert language model at judging whether or not a response adequately addresses an instruction in the context of an 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 which requires a response.
2. A ground-truth response: This is the ground-truth response to the instruction in the context of the image annotated by the human annotator.
3. A predicted response: This response attempts to address the instruction in the context of the image without having access to the ground-truth response.

Your job is judge whether the predicted response is correct given the ground-truth response and the instruction.

Some things to remember:

- Even though you are just a language model, the instructions mostly require an objective answer i.e., the ground-truth response and instruction should be sufficient for you to judge the correctness of the predicted response. You do not need to have access to the complete image description.
- You are capable of judging response quality, accounting for important factors like correctness, relevance, fluency, specificity, etc.
- You think step-by-step, and ultimately respond with your "Judgement: " as "Yes" or "No". Here, "Yes" implies that the predicted response is correct according to you, and "No" implies that the predicted response is not correct.
- Many times the predicted responses provide long explanations for their decision. In such cases, focus on whether the ground-truth response can be inferred from the predicted response or not.

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

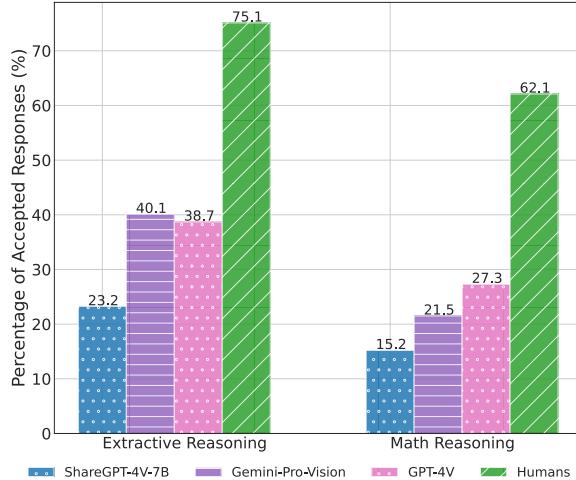
Ground-truth Response: The orangish-brown colored single-seating sofa is the most similar between the Living Room and the Hip Living Room.

Predicted 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.

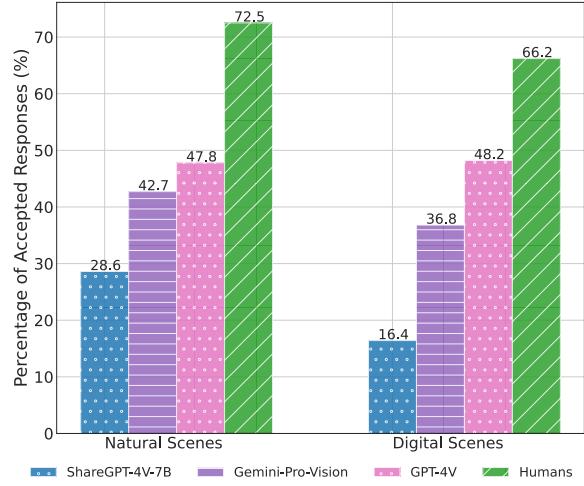
Response: No

Figure 20. Example prompt for GPT4 evaluation. Here, the predicted response is taken from Gemini Pro-Vision

G. Additional Fine-grained Evaluation



(a) Performance on different types of tasks.



(b) Performance on natural and digital scenes.

Figure 21. Additional fine-grained evaluation results.

G.1. Types of Tasks

We compare the performance of the foundation models with varying types of tasks in Figure 21a 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 the large gap in their math reasoning.

G.2. Visual Scenes

We compare the performance of the foundation models with varying visual scenes (e.g., natural scenes and digital scenes) in Figure 21b. 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 the digital scenes. It indicates these models may not have seen many examples with digital scenes during their pretraining. Thus, our CONTEXTUAL dataset highlights the gaps in the training data of the modern LMMs.

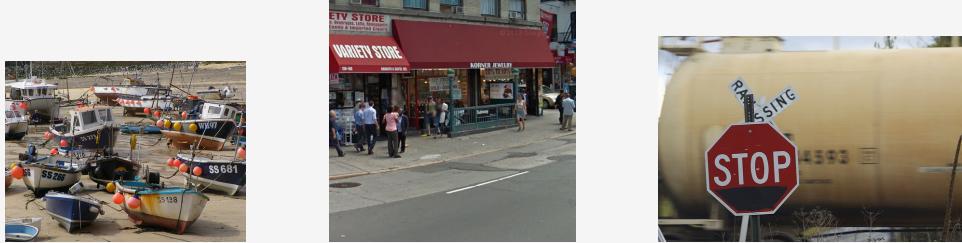
H. Synthetically Scaling Data

We develop an synthetic data generation pipeline that could be useful to create samples that more likely require context-sensitive text-rich visual reasoning. We use the existing OCR filtering strategy to obtain 200 candidate images belonging to Misc. Natural scenes category. Then, we use in-context learning capabilities of GPT-4V and prompt it to generate instruction response pairs (as shown in Fig 22). We observed that use of negative demonstration within in-context examples, as shown in previous studies (Wang et al., 2022), along with reasons was critical for GPT4V to understand the nuanced difference between text only, visual element only and context-sensitive text-rich visual reasoning instructions. After generating the examples, we evaluated representative models from the different model categories: open-source LMMs (LLaVA, ShareGPT), closed-source LMMs (GPT-4V), and Augmented LLM (GPT-4 w/ Layout-aware OCR + Caption), as shown in Table 7.

Model	Synthetic examples	CONTEXTUAL examples
LLaVa-v1.5-13B	39.8	29.7
ShareGPT4v-7B	44.9	37.7
GPT-4 w/ Layout-aware OCR	51.1	27.3
GPT4V	68.6	48.0

Table 7. Comparison of model performance (accuracy in %) using GPT4 evaluation of synthetically generated samples belonging to Misc. Natural Scenes category (200 samples) and human annotated samples (CONTEXTUAL) belonging to Misc. Natural Sciences Category (156 samples).

We observe that the accuracy of models on synthetic samples is greater than the accuracy of human-made ConTextual examples. This is expected because it is difficult for the model to understand context-sensitive in the first place, and asking it to make such tasks would be difficult. Despite this, it can create good instructions because the performance gap of open-source LMMs is relatively small. The only exceptions are Aug LLM and GPT4V, which both use GPT models. Due to this, the GPT4 evaluator may favor their responses to the instruction generated by a GPT model. This demonstrates that we can scale these instructions by filtering good candidates for data generation and nuanced prompt engineering.



Prompt:

You are InstructionResponseGPT4Vision, an expert in understanding images and generating appropriate instruction-response pairs based on exact point Text-Vision Context. The model has access to two images and their corresponding Bad Instruction 1, Bad Instruction 1 reason, Bad Instruction 2, Bad Instruction 2 reason, Good Instruction, Good Instruction response and Good Instruction reason. You need to generate a Good Instruction-Response pair for the third image, that follows the constraints of the Good Instruction. You must always follow the next five guidelines strictly.

1. Instruction must be specific about the details in the image and not be vague or open ended
2. Instructions must always be complex, creative and use different imperative and task oriented main verbs.
3. Note text is not a visual element, a visual element is an object and its attributes / actions/ relationships.
4. An instruction cannot be based on just text as it can be answered using OCR. An instruction cannot be based on just visual elements as it can be answered by vision models. It must always require joint text and visual element reasoning.
5. The response should contain information present in the image, no external or general knowledge can be used or assumptions made.

Input Examples:

Uploaded Image 1:

Bad Instruction 1: Identify the color of the round items on the boats.

Reason: Bad instruction, only requires to color attributes of visual element, here round items, from the image. Does not use joint-text-vision-context to frame the instruction.

Bad Instruction 2: Count the number of red round items if the serial number is SS681.

Reason: Bad instruction, has a conditional "if" in the instruction, conditionals cannot appear in the instruction.

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

Response: WH97

Reason: Good instruction, requires first identifying which boat has three yellow and one round red item and then requires to extract its serial number WH97.

Uploaded Image 2:

Bad Instruction 1: Which stores are visible in this image?

Reason: Bad instruction, as one can obtain the answers using OCR and we do not want that. Does not use joint-text-vision-context to frame the instruction.

Bad Instruction 2: Is the signage of the stores shown in the image of the same color?

Reason: Bad instruction, we only need visual element extraction to respond to the question. Does not use joint-text-vision-context to frame the instruction.

Good Instruction: Name the shop whose main door directly leads to the subway's entrance.

Response: Variety Store

Reason: Good Instruction, it requires to first localize physical elements of the stores, like the main door and see which one is before the subway because that will lead to the subways entrance. After identifying the correct store, it needs to extract the name of that store. This uses joint-text-vision context to frame the instruction.

Now your turn:

Uploaded Image 3:

Good Instruction:

Response:

Good Instruction: Determine the message conveyed by the combination of the STOP sign and the movement implied by the blurred background.

Response: The image conveys the message of a railroad crossing where vehicles are required to stop for passing trains, indicated by the combination of a STOP sign and the motion blur of a train in the background.

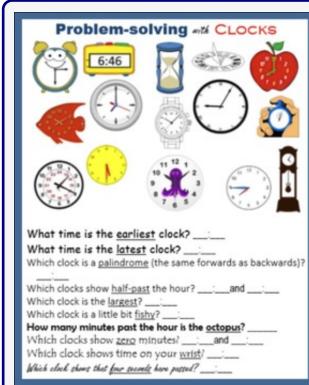
Figure 22. Example prompt for synthetically generating context-sensitive instructions using GPT4V, with images from left-to-right: Uploaded Image 1, Uploaded Image 2, Uploaded Image 3.

I. Detailed Qualitative Analysis

In this section, we qualitatively analyze the responses of the different visual categories. In the caption annotations, we follow the following color notations. **green**: correct response, **red**: incorrect response, and **purple**: intuition behind incorrect response.

I.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.



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**

LLaVA-Next-34B: The object on the dial that shows the time 10:10 is the **hour hand and the minute hand**.

Lack of context-sensitive text-rich visual reasoning

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

Lack of context-sensitive text-rich visual reasoning

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

Lack of context-sensitive text-rich visual reasoning

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.**

Lack of context-sensitive text-rich visual 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-rich visual reasoning. Open source LMMs like *LLaVA-Next-34B*, *ShareGPT4V* and *LLaVA-1.5-13B* and *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) predict the wrong response.



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**.

LLaVA-Next-34B: The sundial in the image indicates that it is approximately **12:30**. The shadow of the gnomon (the central rod) is pointing towards the number 12, which represents noon, and the shadow is about halfway between the numbers 12 and 1, suggesting that it is about 30 minutes past noon. Sundials are used to tell time by the position of the shadow cast by the gnomon, which is aligned with the Earth's axis and the position of the sun.

Reasoning correct but needs accurate fine-grained perception

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

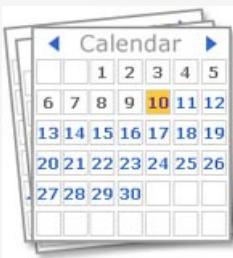
Wrong response

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.**

Lack of context-sensitive text-rich visual 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-rich visual reasoning. *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-rich visual reasoning in this case. Notably, *LLaVA-Next-34B* gets the answer wrong but the reasoning correctly justifies the incorrect answer, thereby highlighting the need for more accurate fine-grained perception.. On the other, *ShareGPT4V* and *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) predict the wrong response.



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.

Partially correct response

LLaVA-Next-34B: The color scheme used for the numbers in the calendar likely serves to highlight specific dates that are important or relevant to the user. For example, the yellow number 10 is highlighted, which could indicate that it is a date of particular significance, such as a birthday, anniversary, or an event that the user wants to easily identify. The use of color in this way helps to quickly draw attention to the important dates, making it easier for the user to remember or plan around them.

Partially 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.

Partially 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-Next-34B & 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.

Hallucination and Incorrect 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

LLaVA-Next-34B: The hour and minute hands on the clock are currently aligned at approximately 1:15. To determine how much time remains until the hour and minute hands align again, we need to find the next time when the hour hand will be at the same position as it is now, which is at the number 1. Since the hour hand moves one number every hour, we can calculate the time remaining by subtracting the current time from the next time the hour hand will be at 1. The next time the hour hand will be at 1 is 2 hours from now, which is 3:15. Now, we subtract the current time from 3:15 to find the time remaining: $3:15 - 1:15 = 2$ hours. So, there are 2 hours remaining until the hour and minute hands align again.

Hallucination and Incorrect 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, *Gemini-Vision-Pro* fails to ground the task to this image. Whereas, *GPT4V*, *LLaVA-Next-34B*, *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.

I.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..



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-rich visual reasoning

LLaVA-Next-34B: The cost of hiring the machine type depicted in the image for nine months is 4500 rupees for a top load machine and **7200 rupees for a front load machine**.

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

Lack of context-sensitive text-rich visual reasoning

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

Lack of context-sensitive text-rich visual reasoning & 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)**

Lack of context-sensitive text-rich visual reasoning

Figure 27. In this task, apart from GPT4V and LLaVA-Next-34B, all other models produce the wrong response. This can be attributed to the strong fine-grained visual perception and context-sensitive text-rich visual reasoning abilities of GPT4V. This also shows that improvement in visual encoding capabilities of open-source LMMs, like LLaVA-Next-34B over LLaVA-1.5-13B improves performance on context-sensitive text-rich visual reasoning tasks.

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.

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.

LLaVA-Next-34B: The most similar product between the living room and the hip living room is **the sofa**. Both images feature a sofa as a central piece of furniture in the living room setup. The color of the sofa in the living room is a neutral gray, while the hip living room has a blue sofa. The other items in the living room and hip living room are different, with the living room having a wooden coffee table and a floor lamp. The hip living room has a more contemporary coffee table and a floor lamp with a unique design.
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.
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 context-sensitive text-rich visual 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-rich visual reasoning 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 context-sensitive text-rich visual 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.

LLaVA-Next-34B: The image shows a group of tourists, but it does not provide an exact number of people. To calculate the maximum money the tour guide can make from the tourists shown, I would need to know the total number of people in the group. If you can provide that information, I can help you with the calculation.

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.

Lack of context-sensitive text-rich visual reasoning

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

LLaVA-Next-34B: The image shows a person holding a camera vertically with a single hand. This type of camera is likely a mirrorless or compact camera, as it appears to be relatively small and lightweight, allowing for such a grip. The camera is equipped with a lens, which is typical for cameras that capture images with a higher level of detail and quality compared to smartphones...

Lack of context-sensitive text-rich visual reasoning

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 ...

Lack of context-sensitive text-rich visual reasoning & Hallucination

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

Lack of context-sensitive text-rich visual reasoning

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...

Lack of context-sensitive text-rich visual reasoning

Figure 30. In this task, all models except *Gemini-Vision-Pro*, show lack of context-sensitive text-rich visual reasoning.

I.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.

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 hallucination.

LLaVA-Next-34B: The triangular-shaped signs in the image are: "LEFT CURVE" in white, "RIGHT CURVE" in white, "LEFT TURN" in white, "RIGHT TURN" in white, "NO ENTRY" in white.
Partially correct response and hallucination.

ShareGPT-4V-7B Response: The triangular-shaped signs in the image are the "No Parking" sign and the "No Entry" sign..
Lack of fine-grained visual perception.

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
Lack of context-sensitive text-rich visual reasoning

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
Lack of context-sensitive text-rich visual reasoning

Figure 31. In this task GPT4V produces the correct response, Gemini-Vision-Pro & LLaVA-Next-34B produce a partially correct answer with minor hallucinations. All other models hallucinate, and produce the wrong response 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.

LLaVA-Next-34B: While driving on the left side of the highway, as indicated by the signage in the image, you can take the following destinations and exits: 1. Exit 24C: South 93 to Quincy 2. Exit 24A: North 93 to Concord, NH 3. Exit 24B: South Station 4. Exit 24A-B: S. Station 5. Clearance 13'-9'6. East 90 7. South Boston 8. Logan Airport. Please note that the specific lanes for each exit may vary, and it's always best to check the road signs and markings for the most accurate information while driving.

Lack of context-sensitive text-rich visual reasoning

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.

Lack of context-sensitive text-rich visual reasoning & Hallucination

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.

Lack of context-sensitive text-rich visual reasoning & Hallucination

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.

Lack of context-sensitive text-rich visual reasoning

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-rich visual reasoning. All other models lack context-sensitive text-rich visual reasoning and some even contain hallucinations in their responses.

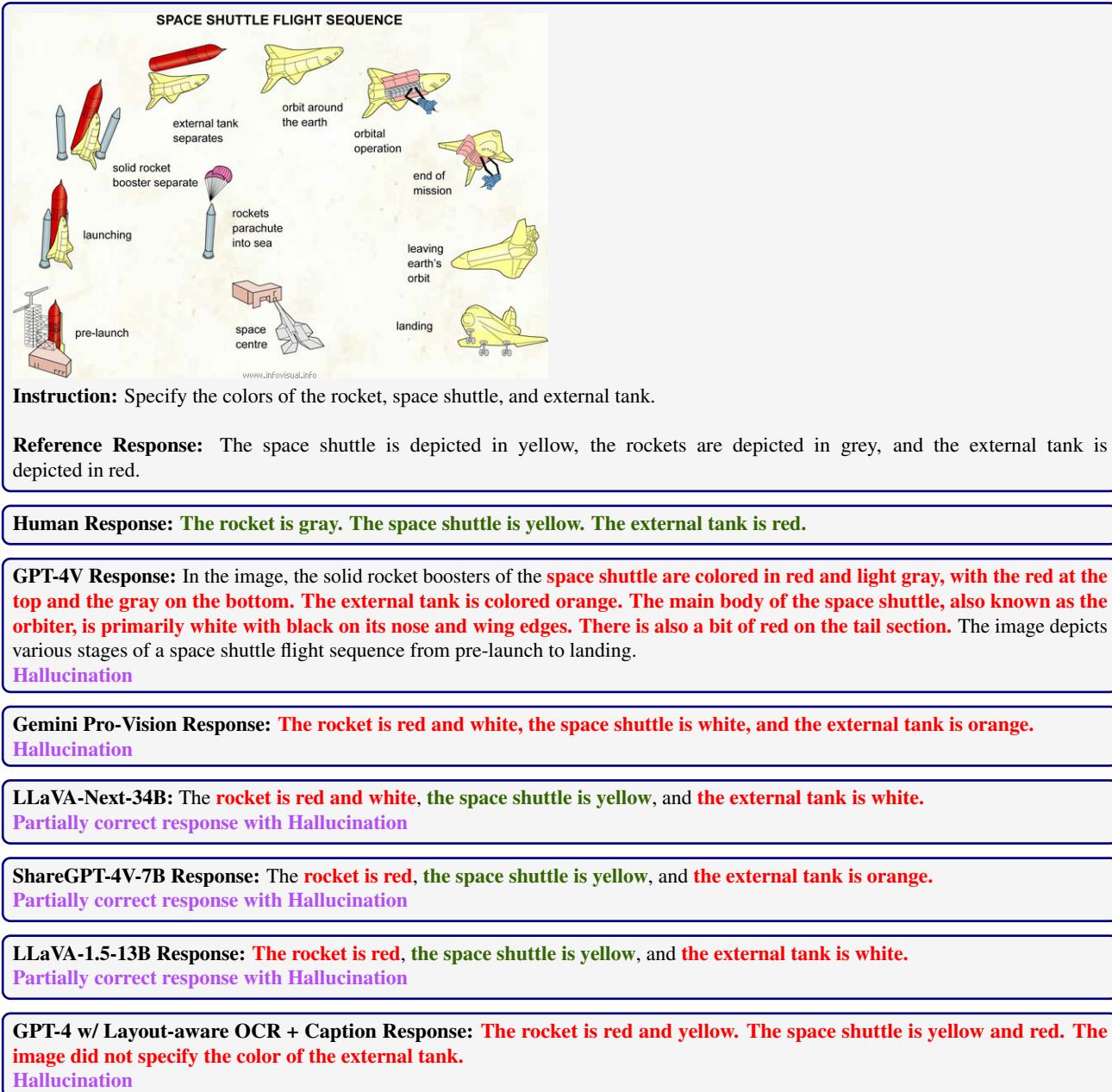


Figure 33. In this task, all models hallucinate, due to lack of grounding the response to the given image & fine-grained visual perception and using prior visual knowledge. Notably, the open-source LMMs, LLaVA-Next-34B, LLaVA-1.5-13B & ShareGPT-4V-7B provide partially accurate responses, suggesting a 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 & Hallucination

LLaVA-Next-34B: The airline immediately in front of the one featuring a flying bird in yellow and blue is "SriLankan". The airline immediately behind it is "Condor"

Lack of fine-grained visual perception & Hallucination

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 & Hallucination

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 & Hallucination

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

I.4. Abstract

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



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.

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.

LLaVA-Next-34B: The image shows a sequence of figures, each depicted in a different pose, with the last figure holding a large ball. The sequence appears to represent a progression from a more challenging task to a more efficient or intelligent approach. The saying "Don't Work Hard, Work Intelligently" is superimposed on the image, which suggests that the figures are illustrating the concept of finding smarter ways to accomplish tasks rather than merely exerting effort. The progression from the first figure, who is struggling with a small box, to the last figure, who is effortlessly holding a large ball, implies that working intelligently can lead to greater success with less effort.

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... 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.' ... Hallucination

Figure 35. In this task, closed-source LMMs GPT-4V, Gemini Pro-Vision & LLaVA-Next-34B correctly respond to the instruction. This can be attributed to their ability for visual perception, text recognition, and context-sensitive text-rich visual reasoning. 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) exclaims "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-rich visual 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-rich visual reasoning

LLaVA-Next-34B: The comic strip features a humorous exchange between two characters, one of whom is a whale. The humor arises from the unexpected and absurd nature of the conversation. The first character asks the whale, "Would the defendant tell the court exactly... what kind of whale you are?" This question is funny because it's an unusual and irrelevant inquiry in a court setting. The whale's response, "Well, I'm... uh...", is funny because it's a human-like reaction to being put on the spot, and the whale's hesitation suggests it's unsure how to answer. The final panel, where the whale objects, is a play on the legal term "objection" and adds to the absurdity of the situation, as whales do not participate in legal proceedings...

Lack of context-sensitive text-rich visual 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-rich visual 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-rich visual 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...

Lack of context-sensitive text-rich visual reasoning & Hallucination

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

I.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.

The screenshot shows a mobile application interface. At the top, there's a header bar with icons and the time '8:25'. Below it is a navigation bar with three items: 'Fabulous', a plus sign, and a search icon. A central card titled 'Drink Water' with a water drop icon says 'Do it 3 times this week to succeed'. It has three circular buttons labeled 1, 2, and 3, with the first one filled. Below this is a section titled 'Today' showing a task 'Morning Ritual' completed at '7:30 AM' with a green checkmark. Under 'Someday', there are two tasks: 'Afternoon Ritual' (Set a reminder) and 'Evening Ritual' (Set a reminder). At the bottom is a navigation bar with arrows and a circular button with a bell icon.

Instruction: How many times has the user drank water today?

Reference Response: 1 times since only the option with "1" on it is selected

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

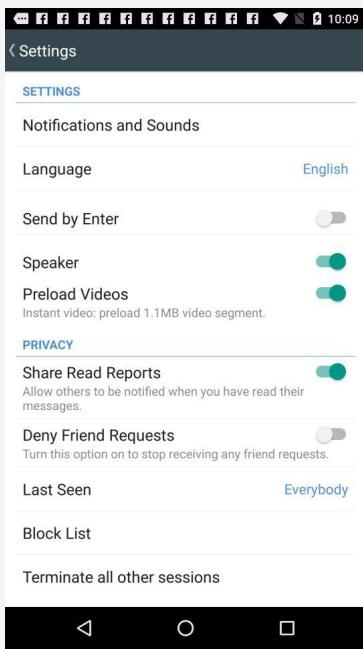
LLaVA-Next-34B: The user has drank water once today.

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* & *LLaVA-Next-34B* produces the correct response, and all other models hallucinate, due to lack of fine-grained visual perception ability. This also shows that improvement in visual encoding capabilities of open-source LMMs, like *LLaVA-Next-34B* over *LLaVA-1.5-13B* improves performance on context-sensitive text-rich visual reasoning tasks.



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

LLaVA-Next-34B: The settings that are switched off by the user are: 1. Notifications and Sounds 2. Language 3. Send by Enter 4. Speaker 5. Preload Videos 6. Share Read Reports 7. Deny Friend Requests 8. Last Seen 9. Block List.

Lack of fine-grained visual perception & Hallucination

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.

Lack of fine-grained visual perception & Hallucination

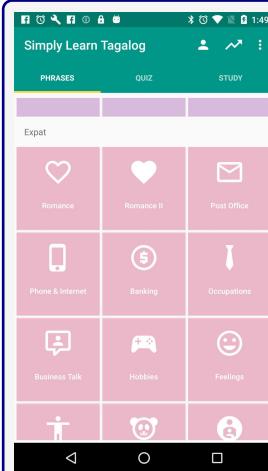
LLaVA-1.5-13B Response: In the image, there is a cell phone with several settings turned off. These settings include the speaker, preloaded videos, and instant video. The speaker setting is turned off, which means that the user may not be able to hear notifications or ringtones. The preloaded videos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

Lack of fine-grained visual perception & Hallucination

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-Next-34B, LLava-1.5-13B & ShareGPT4V-7B suffer from hallucination. 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 visual or thematic change.
Lack of fine-grained visual perception & Hallucination

Gemini Pro-Vision Response: The difference between the two is that the first one has one heart, while the second one has two.
Lack of fine-grained visual perception & Hallucination

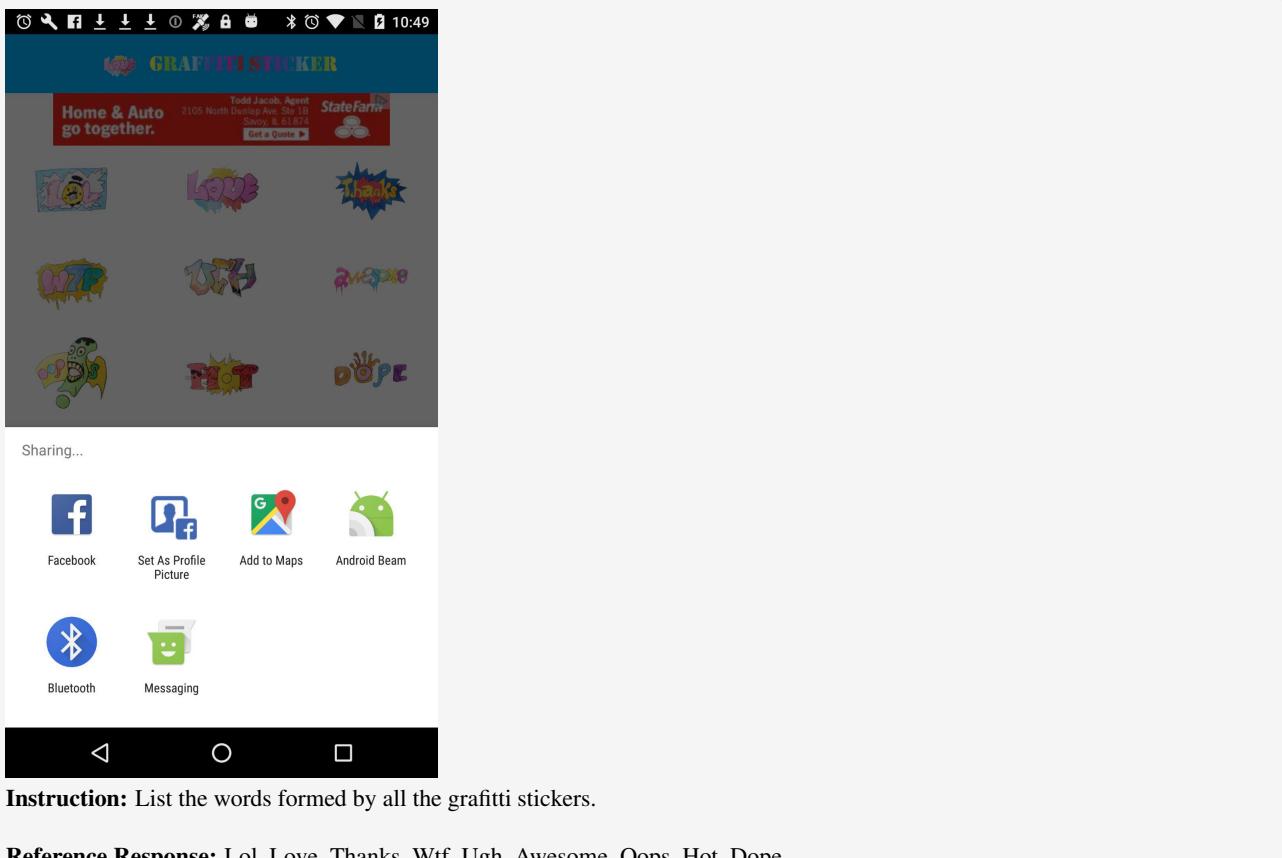
LLaVA-Next-34B: The illustration for "Romance" features a heart symbol, which is a universal symbol for love and affection. The illustration for "Romance II" also features a heart symbol, but it is accompanied by a small lock, which could represent the idea of commitment or the act of locking one's heart to someone. This might suggest that "Romance II" is a more advanced or deeper form of romance, implying a stronger bond or a more serious relationship.
Lack of fine-grained visual perception & 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.
Lack of fine-grained visual perception & 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".
Lack of context-sensitive text-rich visual reasoning

Figure 39. In this task, the closed source LMMs, GPT4V, Gemini-Vision-Pro, LLaVA-Next-34B & LLaVA-1.5-13B hallucinate in their responses. GPT-4 w/ Layout-aware OCR + Caption Response (Augmented LLM approach) produces the wrong answer, due to lack of context-sensitive text-rich visual reasoning. However, ShareGPT4V-7B produces a partially correct response but is unable to ground it completely to the image.



Instruction: List the words formed by all the graffiti stickers.

Reference Response: Lol, Love, Thanks, Wtf, Ugh, Awesome, Oops, Hot, Dope

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.**

Partially correct response with Hallucination

Gemini Pro-Vision Response: **Home, Auto, Wut, Uhh Love,Thanks,Awesome,Dope,Hot**

Partially correct response with Hallucination

LLaVA-Next-34B: The image shows a collection of graffiti stickers, but without a specific word or phrase in mind, it's not possible to list the words formed by them. Graffiti stickers often feature abstract designs, characters, or symbols that may not form words in a traditional sense. If you have a particular word or phrase in mind, please provide more context or details.

Lack of fine-grained visual perception & Lack of context-sensitive text-rich visual reasoning

ShareGPT-4V-7B Response: Love, Awesome, Dope, **Facebook, Google, Android Beam**

Hallucination 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**".

Hallucination

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 text styles.

I.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.

Peter Menzel
Photography

< Books by Peter Menzel & Faith D'Aluisio

Material World: A Global Family Portrait

14 GALLERIES

Material World: A Global Family Portrait, by Peter Menzel, 1994, Sierra Club Books

Search this collection

Material World: Family Portraits 18 IMAGES

Material World: Bhutan 63 IMAGES

Material World: Cuba 8 IMAGES

Material World: Iceland 5 IMAGES

Material World: Japan 37 IMAGES

Material World: Kuwait 3 IMAGES

Material World: Mexico

Material World: Mongolia

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

Lack of fine-grained visual perception

LLaVA-Next-34B: The country whose portrait shows a group of people on the road is Bhutan.

Lack of fine-grained visual perception

ShareGPT-4V-7B Response: Mexico

Lack of fine-grained visual perception

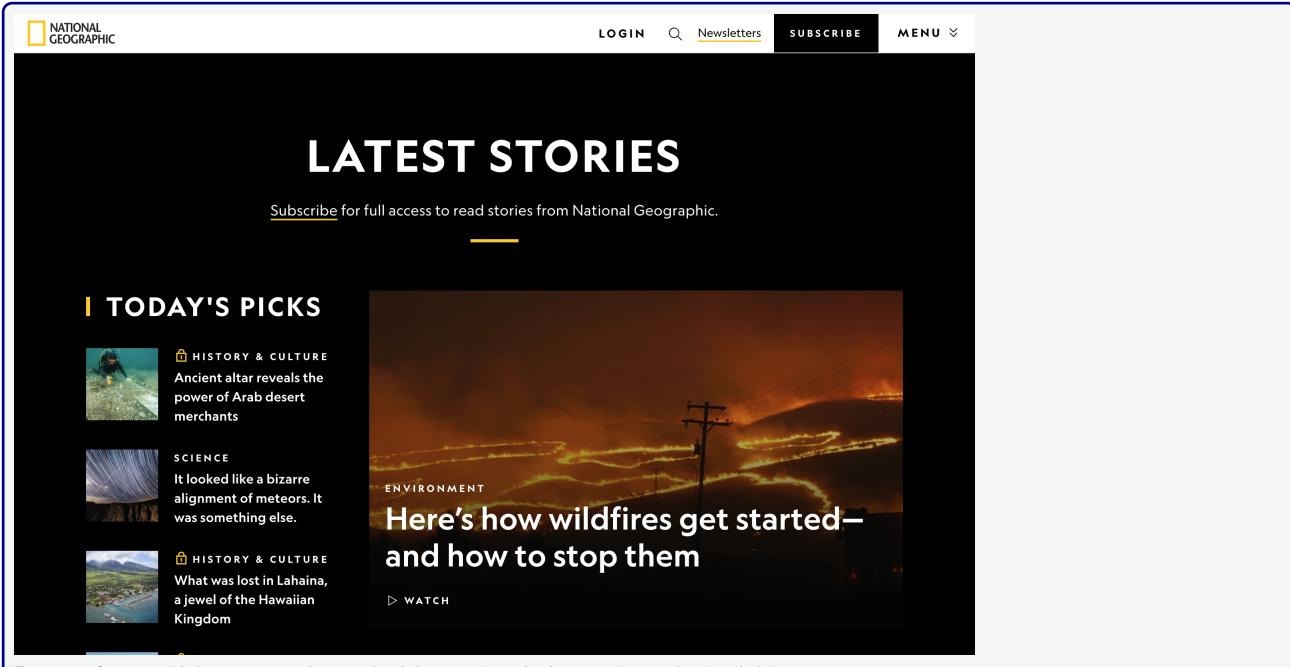
LLaVA-1.5-13B Response: Bhutan

Lack of fine-grained visual perception

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 context-sensitive text-rich visual reasoning

Figure 41. In this task *GPT4V* produces the correct response, and *all other models* produce the wrong answer, due to lack of fine-grained visual perception ability. Specifically, *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) lacks context-sensitive text-rich visual 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.

LLaVA-Next-34B: Under "Today's Picks" on the National Geographic website, there are several stories that can be read without a subscription. These include: 1. "Ancient altar reveals the power of Arab desert merchants" under the "History & Culture" section. 2. "It looked like a bizarre alignment of meteors. It was something else." under the "Science" section. 3. "What was lost in Lahaina, a jewel of the Hawaiian Kingdom" under the "History & Culture" section...

Lack of context-sensitive text-rich visual reasoning

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-rich visual 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-rich visual 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"...

Lack of context-sensitive text-rich visual reasoning

Figure 42. In this task Gemini-Vision-Pro produces the correct response. GPT4V fails to understand the task and ground it to the given image. The open source LMMs, LLaVA-1.5-13B & ShareGPT4V-7B produce the wrong response, due to lack of context-sensitive text-rich visual reasoning. LLaVA-Next-34B & GPT-4 w/ Layout-aware OCR + Caption Response (Augmented LLM approach) hallucinates along with producing the correct response.

Schedule:

8/26 - 11:30 AM PT NBC	8/26 - 2:30 PM PT CBSSN	8/26 - 4:00 PM PT ESPN	8/26 - 4:00 PM PT FS1	8/26 - 4:30 PM PT SECN
NAVY 0-0 ND 0-0	UTEP 0-0 JVST 0-0	MASS 0-0 NMSU 0-0	OHIO 0-0 SDSU 0-0	HAW 0-0 VAN 0-0

College Football Nation Blog:

Georgia's dominant ground game, resurgent D overcome Sooners in 2OT thriller

Scout's Take: ATH Daxton Hill to Michigan

Instruction: Give the jersey number of the leftmost player on the ground from the white team.

Reference Response: 16

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.
Hallucination

LLaVA-Next-34B: The jersey number of the leftmost player on the ground from the white team is 16.

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

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

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

Figure 43. Notably, in this example LLaVA-Next-34B is the only model that gets the correct response. All other models either do not produce a response or produce the wrong response due to lack of strong text recognition capabilities. This also shows that improvement in visual encoding capabilities of open-source LMMs, like LLaVA-Next-34B over LLaVA-1.5-13B improves performance on context-sensitive text-rich visual reasoning tasks.

The screenshot shows a news website interface. At the top, there's a navigation bar with 'AL' logo, 'Alabama', 'Contribute', and 'Sign in'. Below the navigation is an advertisement for 'the Lede'. The ad features a dark background with white text: 'BIRMINGHAM' in a red box at the top left, followed by 'theLede.' in large letters. To the right of the text is a small thumbnail image of two men. Below the thumbnail, the text reads: 'Every story. Every day. Your Birmingham news is on the Lede.' A 'SUBSCRIBE' button is present, with the note 'First month free.' underneath. At the bottom of the page, there's a section titled 'ALABAMA FOOTBALL' with an article about Alabama starting DB dealing with injury after freak accident. The article is attributed to Kevin Scarbinsky and published on Jan. 01, 2018, at 3:01 p.m.

Instruction: Retrieve all the words with white font on a red background.

Reference Response: 1. Birmingham, 2. Subscribe

Human Response: The words with white font on a red background are: 'BIRMINGHAM' and 'SUBSCRIBE'.

GPT-4V Response: The words with white font on a red background are "Every story. Every day." and "First month free." Hallucination

Gemini Pro-Vision Response: Advertisement, theLede, Every story, Every day, First, Your Birmingham news is on the Lede, SUBSCRIBE
Partially correct response with Hallucination

LLaVA-Next-34B: The words with white font on a red background are "Birmingham" and "the Lede"
Partially correct response with Hallucination

ShareGPT-4V-7B Response: Birmingham, Alabama
Partially correct response with Hallucination

LLaVA-1.5-13B Response: Birmingham, Alabama, The Lede, Alabama football, Obama, starting, DB, dealing, with injury, after free president, Alabama and Auburn fans, sign up today.
Partially correct response with Hallucination

GPT-4 w/ Layout-aware OCR + Caption Response: BIRMINGHAM, Advertisement, the, Lede., Every story.Every day, Your, Birmingham, news is, on the Lede., the, Lede, SUBSCRIBE,Birmingham native, Daniel Scheinert, and, All at Once', First month free, ALABAMA FOOTBALL, Alabama starting DB dealing with injury after freak accident, Published: Jan. 01, 2018, 3:01 p.m, Kevin Scarbinsky, The essential newsletter for Alabama and Auburn fans
Partially correct response with Hallucination

Figure 44. In this task, *all models* hallucinate. However, all but *GPT-4V* produce a partially correct response.

I.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.

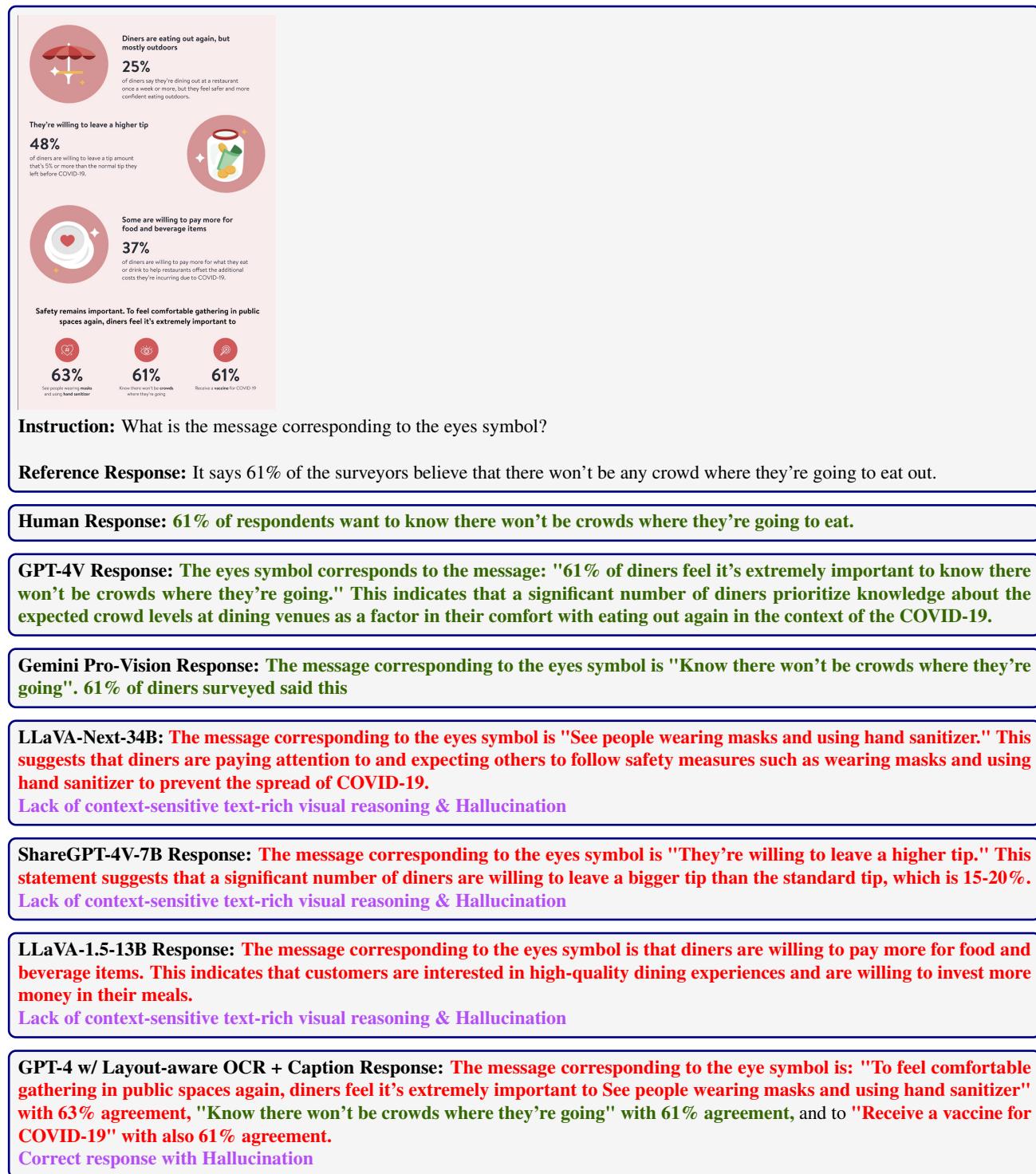


Figure 45. In this task, closed-source LMMs *GPT-4V* and *Gemini Pro-Vision* correctly respond to the instruction. This can be attributed to their ability for visual perception, text recognition, and context-sensitive text-rich visual reasoning. All other models hallucinate.

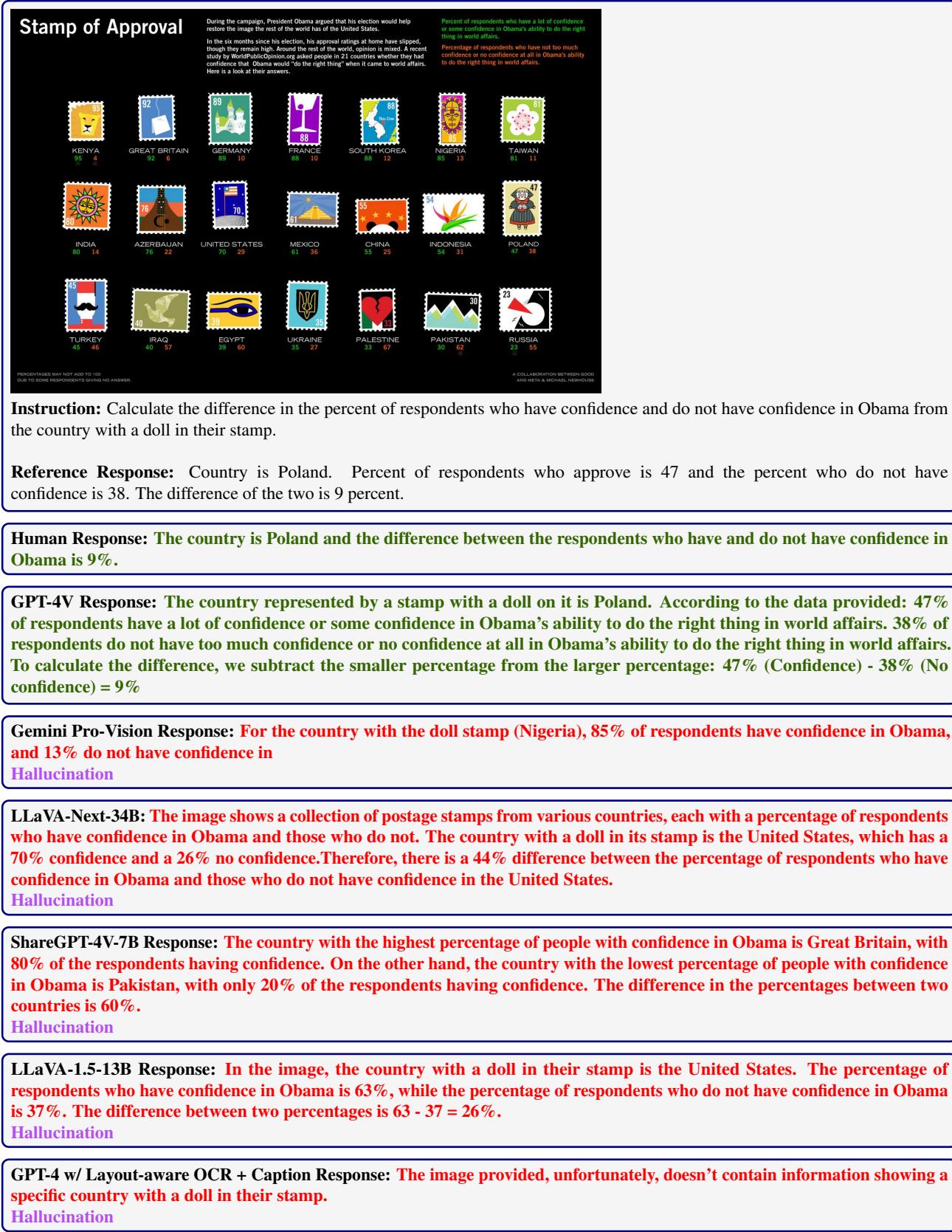


Figure 46. In this task, GPT-4V correctly responds to the instruction. This can be attributed to its ability for math and context-sensitive text-rich visual reasoning. All other models hallucinate their responses.

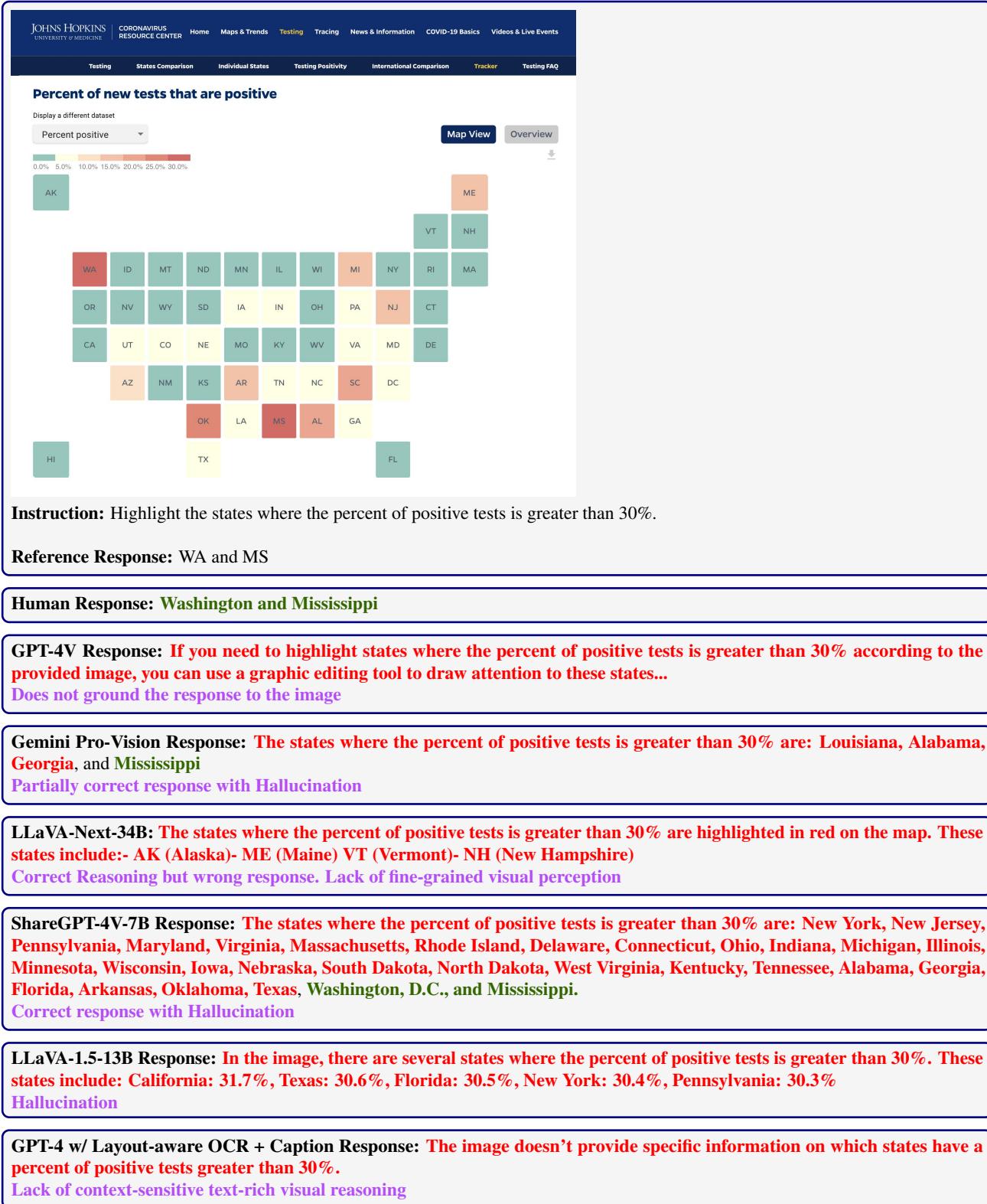


Figure 47. In this task, the closed source LMMs, *GPT4V* does not ground the response to the image. *Gemini-Vision-Pro*, *ShareGPT-4V-7B* & *LLaVA-1.5-13B* hallucinate in their responses. *GPT-4 w/ Layout-aware OCR + Caption Response* (Augmented LLM approach) produces the wrong answer, due to lack of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

Gemini Pro-Vision Response: There are 8 production methods listed under pour-over/drip category.
Lack of fine-grained visual perception

LLaVA-Next-34B: The image shows a chart of coffee production methods, and under the "Pour-Over/Drip" category, there are four different methods illustrated: 1. Chemex 2. V60 3. Kalita Wave 4. Bee House. Each method is represented by a different coffee-making device, and they are connected by lines indicating the flow of the coffee-making process.)
Lack of fine-grained visual perception

ShareGPT-4V-7B Response: There are two production methods used under the pour-over/drip category: French press and drip coffee maker.
Lack of fine-grained visual perception

LLaVA-1.5-13B Response: There are three production methods used under the pour-over/drip category.
Lack of fine-grained visual perception

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 context-sensitive text-rich visual 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, due to lack of context-sensitive text-rich visual reasoning.

I.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.



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.

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

Lack of fine-grained visual perception

LLaVA-Next-34B: Based on the image, the switch is currently in the "OFF" position. If you were to change the switch to the "ON" position, it would likely activate the instrument, assuming that the instrument is in proper working condition and has been set up correctly.

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.

Lack of fine-grained visual perception

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.

Lack of fine-grained visual perception

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 context-sensitive text-rich visual reasoning

Figure 49. In this task, GPT4V & LLaVA-Next-34B 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, due to lack of context-sensitive text-rich visual 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

LLaVA-Next-34B: The image shows a sign with the letters "BM" in large red font, and the word "MUSIC" in smaller white font. The leg of the letter "M" is not visible in the image, so it's not possible to determine what is written on it. If you're referring to the leg of the "M" in "MUSIC," it's not visible in the image provided.

Lack of fine-grained visual perception

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

Lack of fine-grained visual perception

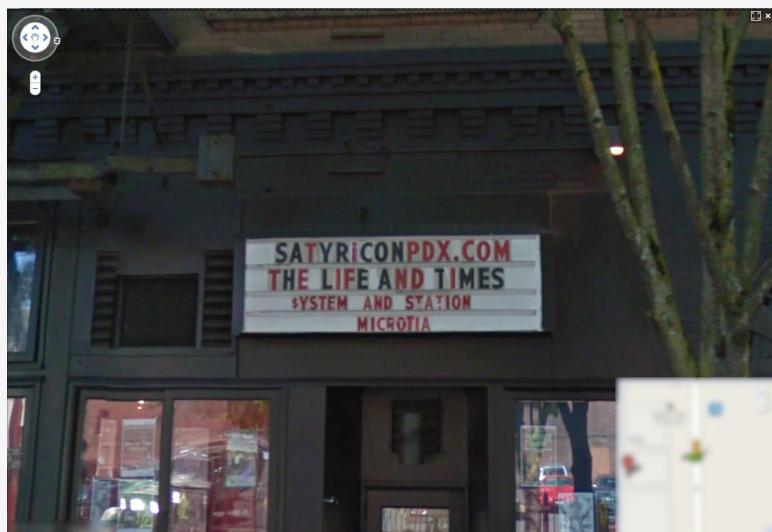
LLaVA-1.5-13B Response: The word "Music" is written on the leg of the letter M.

Lack of fine-grained visual perception

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

Lack of context-sensitive text-rich visual reasoning

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-rich visual reasoning. Open source LMMs like LLaVA-Next-34B, ShareGPT4V & LLaVA-1.5-13B due to lack of fine-grained visual perception. GPT-4 w/ Layout-aware OCR + Caption Response (Augmented LLM approach) predict the wrong answer due to lack of context-sensitive text-rich visual reasoning.



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".

Lack of fine-grained visual perception

Gemini Pro-Vision Response: SATYRICONPDX.COM

Lack of fine-grained visual perception

LLaVA-Next-34B: SATYRICONPX.COM

Lack of fine-grained visual perception

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

Lack of fine-grained visual perception

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

Lack of fine-grained visual perception

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 context-sensitive text-rich visual reasoning

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

Lack of fine-grained visual perception

Gemini Pro-Vision Response: SS 273

Lack of fine-grained visual perception

LLaVA-Next-34B: The boat with three yellow and one red round items hanging from it is the one with the number "SS 681" on it.

Lack of fine-grained visual perception

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

Lack of fine-grained visual perception

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.

Lack of fine-grained visual perception & Hallucination

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 context-sensitive text-rich visual reasoning

Figure 52. In this task, all 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, due to lack of context-sensitive text-rich visual reasoning.