

## 598 Appendix

### 599 A License and Intended Use

600 The VisIT-Bench dataset, along with its various contributions such as instructions, reference outputs,  
601 and model ranking annotations, is licensed under the Creative Commons Attribution 4.0 International  
602 License (CC BY 4.0). This license applies to all the images we have directly contributed, each of  
603 which carries a public license specification in the “public images metadata” field within the dataset  
604 sheets. However, the dataset also incorporates images sourced from pre-existing collections. For  
605 these images, the original licensing terms are respected and remain applicable.

606 VisIT-Bench’s primary purpose is to function as a dynamic benchmark that continuously evolves  
607 and evaluates instruction-following vision-language models. In the current landscape, commercial  
608 chatbots are often trained on non-disclosed and non-public datasets, which raises concerns about  
609 potential data contamination and inadvertent training on our evaluation data [8]. This risk is further  
610 highlighted by recent studies [55, 56]. To mitigate such concerns, we have chosen to withhold the  
611 complete VisIT-Bench test set from public disclosure, while still making the images and instructions  
612 available for direct download. Researchers, however, can utilize VisIT-Bench to its full potential  
613 as a dynamic benchmark by submitting their model predictions for evaluation. We will assess their  
614 models using the undisclosed test set, ensuring the ongoing evolution of the benchmark. Moreover,  
615 we are open to releasing the test data upon receiving reasonable and justified requests, particularly  
616 when additional analysis is necessary, provided that requesters agree to our non-contamination policy  
617 which prohibits the use of this data for training commercial chatbots. This approach strikes a balance  
618 between the need for robust model evaluation and the mitigation of potential data contamination.

### 619 A Dataset Analysis

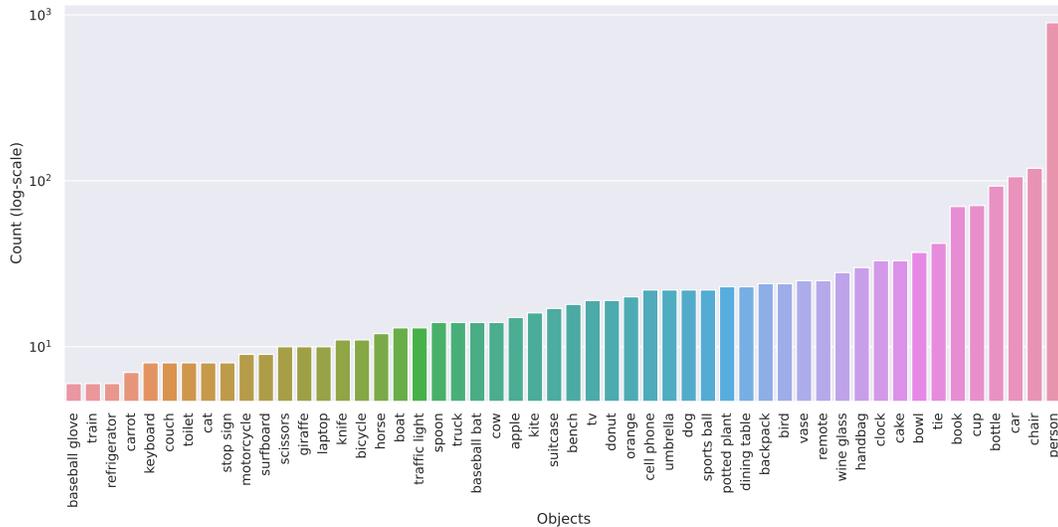


Figure 9: Count of the various COCO objects present in the VisIT-Bench images on a log-scale. The object detection was performed by Yolov5-Large [39]. We observe that the object ‘person’ occurs most of the time in the dataset images i.e., 875 times.

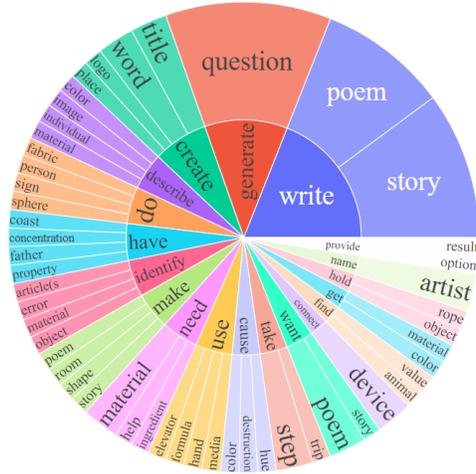


Figure 10: Most frequently occurring verbs (inner circle) and their top 4 direct nouns (outer circle) in the VisIT-Bench instructions.

## 620 B Related Work

621 **Multimodal Models for Image-Text Understanding:** Recently, the field of machine learning  
 622 has experienced a rapid proliferation of new models which can perform various image-text tasks  
 623 [12, 15, 13, 50, 18, 14]. This growth has been driven by several factors, including the emergence of  
 624 large-scale multimodal datasets (e.g. LAION-5B [57], Multimodal C4 [11]), improved software and  
 625 hardware frameworks, and advances in modality-specific models such as language models (e.g., [10]).  
 626 Our work specifically evaluates models which can generate textual outputs, given one or more images,  
 627 and text. Recent examples of such models include LLaVA [13], mPLUG-Owl [17], InstructBLIP,  
 628 LLaMA-Adapter, Flamingo [12] and OpenFlamingo [11], PandaGPT [18], and GPT-4 [7] (which  
 629 reports multimodal capabilities but has not yet seen a release of the multimodal variant).

630 **Instruction Following:** “Instruction-following” is an emerging paradigm for training models via  
 631 language, where instead of being trained to complete only a single, fixed task (such as image  
 632 classification or captioning), models are trained to follow textual instructions that describe an  
 633 arbitrary task, with the aim of generalizing to novel instructions. Examples of instruction-following  
 634 models include Alpaca [5], LLaMA-Adapter [16], Koala [51], InstructBLIP [14], LLaVA [13], and  
 635 mPLUG-owl [17]. As the downstream capabilities of these models are influenced by the quality of  
 636 the training dataset, there has also been extensive work on developing instruction-following datasets  
 637 [38, 58, 59, 13, 60].

638 To build powerful these models, two broad approaches have been shown to be effective. One  
 639 approach focuses on leveraging existing pretrained task-specific tools such as image captioners  
 640 [15], object detectors [61] and text-to-image generators [62] by either creating multimodal prompt  
 641 interfaces [42, 63] or by executing LLM-generated programs [64, 65, 66]. The other approach  
 642 [13, 16, 67, 45, 68, 17, 11] focuses on building a single pretrained model that can follow instructions  
 643 by supervised finetuning on multimodal vision-language data.

644 Despite the success of both these approaches on the existing vision-language datasets e.g., VQA,  
 645 GQA, Image Captioning [21, 52, 20], there is a lack of a high-quality benchmarking dataset for  
 646 multimodal instruction-following tasks that reliably replicates the way in which humans would  
 647 interact with multimodal chatbots in the wild. Similar to the image-text models discussed above,  
 648 many instruction-following models have been released directly as open-source without undergoing  
 649 peer review or thorough evaluation. As a result, the effectiveness of these models for many tasks is  
 650 not well-understood.

651 **Benchmarks for Machine Learning:** High-quality evaluation datasets have served both to (re)assess,  
 652 and to accelerate, progress on many machine learning tasks [53]. For example, our work draws  
 653 particularly from the fields of computer vision and natural language processing, where benchmarking  
 654 datasets have been critical drivers of progress. On the vision side, datasets such as ImageNet [69] and  
 655 CIFAR [70] have proven to be critical yardsticks of progress. On the language side, benchmarks such  
 656 as SQuAD [71], SST [72], GLUE/SuperGLUE [73, 74] and more [75, 76] seen wide use. Recent  
 657 work has indicated that improvements on these high-quality benchmark datasets is not the result of  
 658 overfitting, and is a reliable indicator of genuine progress beyond the benchmark data [77, 78, 79, 80].

659 However, high-quality benchmarking datasets and evaluation methods do not yet exist for multimodal  
 660 instruction-following. As a result, it is difficult to assess progress in this direction, which both reduces  
 661 the field’s ability to identify true breakthroughs and increases vulnerability to potential pitfalls of  
 662 evaluation that have hampered progress in other areas of machine learning [53, 54].

## 663 C Interfaces for Collecting Human Annotations

664 In this section, we provide the templates we used to collect human annotations for the instruction  
 665 generation (Figure 11), the dense caption generation (Figure 12), the model verification (Figure 13  
 and Figure 14), and the model rating (Figure 15).

	<p><b>Instruction</b></p> <p>What is this dish and how can I make it?</p>	<p><b>Expected Output</b></p> <p>This dish appears to be Massaman Curry in a large pot, add onion and sauté, with curry paste, ginger and garlic. 3. Sauté chicken and potatoes are cooked through in 5 minutes. 5. Serve with rice (white, brown) and cilantro on top.</p>
<p><b>Find an Image to Annotate to test Food Recipe, and create a Visual instruction, that a textual-only model/human won't</b></p>		
<p>Feel free to skip the HIT if you don't feel creative with this one</p>		
<p>Provide a <b>**PUBLICALLY LICENSED**</b> image URL.</p>	<p>0/200</p>	
<p>Provide a <b>**creative and challenging**</b> model instruction that would blow your mind if an AI model could answer it.</p>	<p>0/200</p>	
<p style="text-align: center;"> <span>Previous HIT</span>    Showing Task 23 of 70    <span>Next HIT</span> </p>		

Figure 11: An interface that collects the image and corresponding instruction.

	<p><b>Instruction (Tested Skill=Exercise)</b></p> <p>Based on commonly held beliefs about upper body strength in men and women which competitor is likely to win this competition.</p> <p>A highly visual description that will allow TEXT-ONLY models *MULTIPLE* answers for the instructions WITHOUT THE IMAGE</p>
<p>0/2000</p>	

Figure 12: An interface that collects the dense caption given an image-and-instruction pair.

666

**Dense Caption:**  
A woman wearing a red gown is sitting in a wheelchair but dancing with a man in a black tuxedo and white dress shirt. She is using her left hand to hold his right hand while posing with her other hand. He has left arm stretched out. They are performing on a stage with purple highlights.

**Instruction:**  
Why is seeing this couple compete in this event an uncommon sight? Does this convey a positive message?

**GPT-4 Prediction:**  
Seeing this couple compete is an uncommon sight because the woman is in a wheelchair, which is not typical in traditional dance performances. However, this conveys a positive message as it showcases inclusivity, determination, and the ability to overcome challenges.

Does GPT4 correctly follow the instruction?  
 Yes  
 No

If answered 'No' for the above, mark one of the following:  
 1.  N/A  
 2.  The dense caption is detailed enough; the problem is in GPT4's answer  
 3.  The dense caption is underspecified; not GPT4's fault

\* Might be offensive/bad advise/problematic

Submit

Figure 13: A sample of the model verification where the GPT4 follows the instruction correctly.

**Dense Caption:**  
There are two equations on a white paper: 1)  $f(x) = \sqrt{x+6} - 1$  and 2)  $f(x) = \frac{1}{x+2} - 1$

**Instruction:**  
For the following exercises, graph the functions by translating, stretching, and/or compressing a toolkit function.

**GPT-4 Prediction:**  
To graph the functions, follow these steps: 1) For  $f(x) = \sqrt{x+6} - 1$ , start with the toolkit function  $y = \sqrt{x}$ . Translate it 6 units to the left and 1 unit down. 2) For  $f(x) = \frac{1}{x+2} - 1$ , start with the toolkit function  $y = \frac{1}{x}$ . Translate it 2 units to the left and 1 unit down.

Does GPT4 correctly follow the instruction?  
 Yes  
 No

If answered 'No' for the above, mark one of the following:  
 1.  N/A  
 2.  The dense caption is detailed enough; the problem is in GPT4's answer  
 3.  The dense caption is underspecified; not GPT4's fault

\* Might be offensive/bad advise/problematic

Submit

Figure 14: A sample of the model verification where the GPT4 failed to follow the instruction due to its incorrect reasoning.

**Instruction:** What is the person doing looking at the laser ?

**A:** The person is examining the red laser, possibly trying to figure out its purpose or how to avoid it. He might be in a high-security area or facing an obstacle.

**B:** The person is standing in front of a laser with their back to the camera. They are looking at the laser, possibly assessing or examining its features.

**Which model better fully follows the instruction, providing factual, helpful and accurate response, and adhering to all of the instruction details?**

Model A  
 Model B

Submit

Figure 15: An interface that collects the feedback of the model rating.

667 **D Existing Datasets incorporated in VisIT-Bench**

668 In Table 5, we listed the existing datasets that are incorporated in our VisIT-Bench. Among these  
 669 datasets, 15 contain a single image in each sample pair, and 10 require reasoning based on multiple  
 670 images.

Table 5: List of existing datasets in VisIT-Bench, categorized as single and multiple image datasets.

	Dataset	Topic
Single	VQA [21]	Visual Question Answering
	VCR [81]	Cognition-level Visual Understanding
	GD-VCR [82]	Geo-Diverse Commonsense Reasoning
	WHOOOPS [40]	What Makes this Image Strange
	Newyork Caption [83]	Humor Understanding
	CLEVR [84]	Visual Question Answering
	Kilogram [85]	Tangrams Identification
	Harmful Memes [86]	Memes Understanding
	ScienceQA [87]	Science Question Answering
	OK-VQA [88]	Outside Knowledge Visual Question Answering
	AOK-VQA [89]	Outside Knowledge Visual Question
	AOK-VQA [89]	Question Generation
	VizWiz [90]	Visual Question Answering
	GQA [52]	Visual Question Answering on Scene Graphs
	TextCaps [91]	Visual Question Answering on Texts
Multiple	Robust Change Captioning [22]	Describing What has Change in a Scene
	NLVR2 [36]	Testing Visual Language Bias
	ImageCoDE [92]	Image Retrieval
	Spot-the-Diff [93]	Identifying Differences
	VASR [94]	Visual Analogies
	WinoGavil [95]	Visual Associations
	IRFL (Metaphor) [96]	Figurative Speech Understanding
	IRFL (Idioms) [96]	Figurative Speech Understanding
	IconQA [97]	Abstract Diagram Understanding
	Pick-a-Pic[98]	Text-to-Image User Preferences

Table 6: List of skills and existing datasets in VisIT-Bench

‘scienceqa’, ‘ocr math’, ‘recognition’, ‘okvqa’, ‘house plan understanding’, ‘nlvr2’, ‘gardening tips’, ‘textcaps’, ‘architectural styles’, ‘dressing sense’, ‘winoground’, ‘food recipe’, ‘paper folding’, ‘whoops’, ‘spot the diff’, ‘wino- gavil’, ‘imagecode’, ‘exercise’, ‘art knowledge’, ‘gqa’, ‘physical knowledge’, ‘contextual knowledge of events’, ‘home renovation’, ‘aokvqa’, ‘animals’, ‘vasr’, ‘counting’, ‘board games’, ‘solving geometry problems’, ‘who to call?’, ‘clevr’, ‘building materials’, ‘hazard identification’, ‘pickapick’, ‘astronomy’, ‘figurative speech explanation’, ‘write a story’, ‘gestures understanding’, ‘newyork’, ‘cul- tural knowledge’, ‘aokvqg’, ‘traffic sign identification’, ‘pop culture’, ‘fashion products’, ‘harmful memes’, ‘write a poem’, ‘vizwiz’, ‘guesstimate of capacity’, ‘location understanding’, ‘graph reasoning’, ‘vqa’, ‘game playing’, ‘differently abled’, ‘chemical identification’, ‘history knowledge’, ‘climate and weather understanding’, ‘irfl metaphor’, ‘human emotion recognition’, ‘medical’, ‘gd vcr’, ‘vcr’, ‘technical support’, ‘catchy titles’, ‘kilogram’, ‘anagrams’, ‘color’, ‘tour guide’, ‘directions’, ‘irfl idiom’, ‘rcc’
--

## 671 **E Elo Rating**

672 For many years, the Elo rating has been popular in ranking players in zero-sum games such as chess  
673 [25]. Recently, it has been adopted to rate large language models (LLMs) against each other on the  
674 user instructions. In this work, we adopt the same strategy to rank a set of instruction-following  
675 vision-language models, that can grow dynamically with further advances in the field.

676 Given two multimodal chatbots  $C_a$  and  $C_b$  with their absolute Elo rating  $\mathcal{R}_a$  and  $\mathcal{R}_b$ , respectively.  
677 Simply put, the probability of  $C_a$  winning over  $C_b$  in a head-to-head battle is given by:

$$P(C_a \text{ wins over } C_b) = \frac{1}{1 + 10^{(\mathcal{R}_a - \mathcal{R}_b)/400}} \quad (1)$$

678 In practice, calculating the Elo rating requires us to set hyperparameters to decide the weightage  
679 for each win and loss in a head-to-head battle between two models. In our work, we use the open  
680 implementation of Elo for LLMs by FastChat at [https://github.com/lm-sys/FastChat/blob/  
681 main/fastchat/serve/monitor/elo\\_analysis.py](https://github.com/lm-sys/FastChat/blob/main/fastchat/serve/monitor/elo_analysis.py).

## 682 **F GPT-4 Pairwise Evaluation Prompts**

683 The specific prompts we use to extract pairwise judgements from our language model are provided in  
684 Table 16 (reference-free version) and Table 17 (reference-backed version). When applied to GPT-4  
685 [7], these prompts usually solicit a definitive pairwise response by the model. But, in some cases, the  
686 model either produces a pairwise judgement in an unexpected format, or, refuses to issue a judgement  
687 at all. For cases like these, we issue an additional query to ChatGPT to extract an answer (or decide  
688 there is no answer) using an additional prompt, given in Table 18. If after this step there is still no  
689 definitive pairwise judgment, we call the result a tie.

690 [UGFt]A head-to-head comparison involves querying GPT-4 two times. Each query has approxi-  
691 mately 800 input/250 output tokens (output includes chain of thought; see Supplementary for full  
692 prompt/response example). Running a head-to-head comparison between a new model and a ran-  
693 domly sampled existing model across the 700 instances is usually sufficient for a stable Elo estimate.  
694 In total, at current GPT-4 prices, the 700 head-to-head comparisons required to assess a new model  
695 costs 50-70 dollars.

---

system prompt (human authored)

You are ImageTaskEvaluationGPT, 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 image context: This will describe the contents of an image with sufficient detail to address the instruction.
2. An instruction: This is a question, an imperative request, or something similar about the image which requires a response.
3. Two responses, response A and response B: These two responses attempt to address the instruction in the context of the image. Your job is to judge whether response A or response B better. A and B are randomly ordered.

Some things to remember:

- Even though you are just a language model, the image description will be sufficiently detailed so that your judgements can be accurate.
  - You are capable of judging response quality, accounting for important factors like correctness, relevance, fluency, specificity, etc.
  - You think step-by-step, but ultimately respond with "Response A" or "Response B"
- 

user (human authored)

I will describe an image to you, and provide an instruction. Then, I will give you two candidate responses that address the instruction in the context of the image: these will be labelled "Response A" and "Response B". Your job is to first reason step-by-step about which response is best in terms of accuracy, specificity, fluency, etc. After reasoning step-by-step and comparing the pros/cons of each response, in the end, respond with "Overall, Response X is better." where X is either A or B.

assistant (human authored)

Sure, please provide the image context, the instruction, and the two candidate responses, Response A and Response B. Then, I will think step-by-step and provide my ultimate judgement as to which response is better.

user (from human-authored template)

OK. Here is the image description, the instruction, and the two response options, Response A and Response B. Image context: A skier is wearing a green jacket, burgundy pants, white shoes, white gloves, and a gray beanie. There are people sitting on the background with their snowboards.

Instruction: Answer the following question based on the provided image: Are the trousers the same color as the shoes?

Response A: no

Response B: No, the trousers are not the same color as the shoes. The trousers are burgundy, while the shoes are white.

Think step-by-step and finish your response with "Overall, Response X is better." where X is either A or B.

assistant (machine generated; used to derive the prediction)

Step 1: Identify the key information from the image description.

The skier is wearing:

- Burgundy pants
- White shoes

Step 2: Assess the accuracy of each response in addressing the instruction.

Both Response A and Response B convey that the trousers and shoes are not of the same color.

Step 3: Compare specificity and detail of each response.

Response A is concise and provides a straightforward answer, but lacks detail.

Response B provides a clear, comprehensive answer that includes the color of both the trousers and shoes.

Step 4: Consider the fluency and correctness.

Both responses are fluent and grammatically correct.

Overall, Response B is better.

---

Figure 16: The prompt we use for reference-free evaluation, along with a sample instance/completion from GPT-4 [7].

---

system prompt (human authored)

You are ImageTaskEvaluationGPT, 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 image context: This will describe the contents of an image with sufficient detail to address the instruction.
2. An instruction: This is a question, an imparative request, or something similar about the image which requires a response.
3. A reference output: This is a high-quality example output that humans have judged to be an accurate response for the input instruction.
4. Two responses, response A and response B: These two responses attempt to address the instruction in the context of the image.

Your job is to judge whether response A or response B better. A and B are randomly ordered. Some things to remember:

- Even though you are just a language model, the image description will be sufficiently detailed so that your judgements can be accurate.
- Take the high-quality reference into account when making your judgements, but remember: some instructions are more open-ended than others, so for those cases, a high quality response can differ from the reference.
- You are capable of judging response quality, accounting for important factors like correctness, relevance, fluency, specificity, etc.
- You think step-by-step, but ultimately respond with "Response A" or "Response B"

---

user (human authored)

I will describe an image to you, and provide an instruction. Then, I will provide a reference output which is an example of a high quality output for that instruction in the context of the image. Then, I will give you two candidate responses that address the instruction in the context of the image: these will be labelled "Response A" and "Response B". Your job is to first reason step-by-step about which response is best in terms of accuracy, specificity, fluency, etc. After reasoning step-by-step and comparing the pros/cons of each response, in the end, respond with "Overall, Response X is better." where X is either A or B.

assistant (human authored)

Sure, please provide the image context, the instruction, the reference, and the two candidate responses, Response A and Response B. Then, I will think step-by-step and provide my ultimate judgement as to which response is better.

user (from human-authored template)

OK. Here is the image description, the instruction, the high-quality reference, and the two response options, Response A and Response B.

Image context: There is a white plate that has on it a small pizza with many toppings, mostly unrecognizable, due to the pie being very well done, charred to the point of almost being burnt, and the crust blackened in many places.

Instruction: Someone told the creator of this item that they are very talented. What is the implied meaning?

High quality reference: The implied meaning is that the person complimenting the creator is being sarcastic, as the pizza is overcooked and almost burnt, which is not a sign of talent.

Response A: When someone says that they are "very talented" in the context of the image, it implies that the person responsible for creating the pizza is skilled in the craft and able to create an appealing and tasty dish. The person praising the pizza's creator recognizes their culinary skills, which can be seen through the combination of ingredients and the overall look of the finished dish. The praise is a form of appreciation for their efforts and the enjoyable result they achieved.

Response B: It is difficult to determine the implied meaning without further context or information about the image. It appears to be a pizza that has been burnt or charred on the surface and has various toppings, including pepperoni, onion, and mushrooms. The image is shown on a white plate or surface. The implied meaning of this image could be that it is a dish that has been cooked badly or that it is a dish with a unique or unusual combination of ingredients.

Think step-by-step and finish your response with "Overall, Response X is better." where X is either A or B.

assistant (machine generated; used to derive the prediction)

Step 1: Comparing both responses with the high-quality reference, it's clear that Response A describes the pizza as well-made and appealing, while Response B is more focused on the fact that it is burnt or charred, somewhat aligning it with the reference.

Step 2: Evaluating specificity, Response A defines the creator as skilled, even though the image description contradicts this by stating that the pizza is overcooked. Response B, while not as specific, acknowledges the overcooked nature of the pizza.

Step 3: In terms of fluency, both responses are easily understood and written in coherent language.

Step 4: Considering accuracy, Response A is inaccurate in relation to the context and instruction because it implies a positive meaning when it should be sarcastic. Response B doesn't fully capture the sarcasm but is closer to the truth about the pizza being badly cooked.

Step 5: Understanding relevance, Response A is not entirely relevant due to the inaccuracy in the creator's talent, while Response B, although less direct, is more relevant by mentioning the burnt or charred nature of the pizza.

Overall, Response B is better.

---

Figure 17: The prompt we use for reference-backed evaluation, along with a sample instance/completion from GPT-4 [7].

---

system prompt (human authored)

You are FinalAnswerExtractionGPT, an expert language model at extracting multiple choice answers from written out from longer explanations. You will be given several sentences describing a thought process which should eventually contain a final answer, either A or B. Your job is to extract the final answer that the explanation arrives at.

Some things to remember:

- Keep your answers short: only output "Final Answer: X" where X is A or B
- If the explanation doesn't mention a final choice explicitly, you can output "Unknown"

...

(+ 2 in-context examples (human authored))

...

---

[Example One: GPT-4 gives a definitive answer, but in the wrong format](#)

user (GPT-4's raw pairwise judgement)

Step 1: Compare the two responses to the instruction

- The instruction asks if people are required to stop and requests an explanation.

Step 2: Assess Response A

- Response A indicates that people should stop due to the presence of a stop sign. However, it doesn't consider the fact that the sign is pressed against the bus, and the LED lights are off.

Step 3: Assess Response B

- Response B simply says "yes," acknowledging that people should stop, but it doesn't provide the explanation required by the instruction.

Step 4: Compare specificity, fluency, and other factors

- Response A is more specific and fluent but misconstrues the image description.
- Response B is accurate in that it acknowledges the stop sign but lacks an explanation.

Overall, since neither response is fully correct, it is difficult to declare a better response. However, given the fact that Response B acknowledges the prominent stop sign, I would say that Response B is slightly better, but it should also explain how the situation indicates that people should stop.

Please extract the final answer from the above text.

assistant (machine generated; used to derive the prediction)

Final Answer: Response B is slightly better, but it should also explain how the situation indicates that people should stop.

---

[Example Two: GPT-4 refuses to answer.](#)

user (GPT-4's raw pairwise judgement)

Step 1: Consider the image context, which states that the original thirteen colonies are shown with New Hampshire specifically highlighted in a darker green.

Step 2: Assess the accuracy of Response A, which claims that the colony shown is "Virginia." This is not accurate since the image context indicates New Hampshire being highlighted.

Step 3: Assess the accuracy of Response B, which claims that the colony shown is the Massachusetts Bay Colony. This response is also not accurate, as the image context states New Hampshire is the highlighted colony.

Step 4: In terms of accuracy, both responses are incorrect. Neither response is consistent with the image context.

Step 5: Consider other factors like fluency and specificity. Both responses seem fluent and specific in nature.

After considering these factors, none of the responses are accurate given the information provided in the image context. However, I am unable to provide an "Overall, Response X is better." judgement in this case as both responses are incorrect and neither stands out as being better than the other.

Please extract the final answer from the above text.

assistant (machine generated; used to derive the prediction)

Final Answer: Unknown.

---

Figure 18: The prompt we use to extract answers from ill-formatted GPT-4 [7] responses, along with two example outputs. In one case, GPT-4's pairwise judgement is given, but in a format different than requested. In the other case, GPT-4 "defies" the prompt, and refuses to issue a judgement because both options are bad.