PhD: A ChatGPT-Prompted Visual hallucination Evaluation Dataset

Supplementary Material

10

12

14

14

In the supplementary materials, we report

- Manual assessment of the dataset quality (Sec. S1);
- ChatGPT instructions for dataset construction (Sec. S2);
- Additional experimental results (Sec. \$3);
- PhD question-answer data (data.json) (Sec. S4).

S1. Manual Assessment of Dataset Quality

We assessed the quality of PhD by randomly selecting 100 samples (*i.e. image-question-answer* triplets) per mode / task. Each sample has a specific image associated with a Yes question and a No question. As an image might have multiple questions while a sample is shared across modes and tasks, the subset has 885 distinct images and 1,794 samples in total, see Tab. S1. Manual assessment was conducted with the following evaluation criteria:

- The VQA data must be correct, with a correct Yes / No label w.r.t. the given image and question;
- The question (and the preceded context) shall effectively induce hallucinations, with the hitem, context, and CCS images designed to be challenging.

Samples meeting both criteria are labeled as good, and bad otherwise. Our annotation team was comprised of 9 volunteers (4 phd and 5 graduate students in our lab), who have good knowledge about computer vision and deep learning, yet do not involve in the current project.

The labeling result is analyzed and summarized in Tab. S1, showing that the sampled subset of PhD achieves a good rate of 93%. Note that the annotation process has been conducted in a collective manner where a specific sample was assessed independently by different annotators, labeling discrepancy naturally exists. The relatively small deviation of 2% on average w.r.t. the good rate indicates a consensus of opinions among the annotators.

	#images	#VQA Triplets	#good	good-rate
Mode				
PhD-base	100	200	188	$94\%{\pm}2\%$
PhD-sec	100	200	185	93%±3%
PhD-icc	100	200	192	96%±2%
PhD-ccs	100	200	192	96%±1%
Task				
Object recog.	100	200	189	95%±2%
Attribute recog.	100	200	188	$94\%{\pm}2\%$
Sentiment recog.	100	200	179	90%±3%
Positional recog.	100	200	178	89%±4%
Counting	100	200	184	$92\%{\pm}3\%$
Overall	885	1,794	1,675	93%±2%

Table S1. **Dataset quality assessment**. #good indicates the number of sampled triplets manually assessed as good.

S2. ChatGPT Instructions for Dataset Construction

ChatGPT instructions used in each step of the proposed pipeline are listed as follows.

S2.1. Instructions for Vocabulary Construction

```
Targeted word dictionary is like:
   "object": ['car', 'bus', 'person' ...],
   "attribute_color": ['red', 'blue',
       · · · ],
   "attribute_material": ['plastic', ...],
   "attribute_shape": ['square', ...],
   "sentiment": ['happy', ...]
   # "counting": [...] zero to nine
   # "position": [...] similar to object
 }
Please expand the input dictionary as
    much as possible by adding common
    items found in daily life under each
    key. Avoid any duplication.
 Output the result directly in <Python
    dictionary> format without providing
    any explanation or comments.
```

output:

Prompt S1. Vocabulary construction

S2.2. Instructions for Subject-Attribute Extraction from TDIUC Annotations

Based on the provided question <q> and < task>, extract the specific subject</q>
being asked about in the question. Output requirement: please output only the subject without any explanation or additional information.
Here are some examples: task: positional_reasoning q: What is to the right of the book? output: to the right of the book
<pre>task: attribute_color q: What is the color of the grass? output: grass</pre>
task: attribute_material q: What is the chair made of?

```
output: chair
                                                       output: Is there a person to the left of
15
                                                    28
                                                           the TV in the image?
16
  task: sentiment_understanding
                                                    29
   q: Is this man happy?
                                                       input: happy
18
                                                    30
  output: man
                                                       related object: dog
                                                    31
19
                                                       task: sentiment_understanding
                                                    32
20
  task: counting
                                                    33
                                                       output: Is the dog happy in the image?
21
  q: How many people are there?
                                                    34
22
  output: people
                                                       input: 3
                                                    35
                                                       related object: bags
24
                                                    36
  Now, following the examples above, please
                                                       task: counting
25
                                                    37
       extract the relevant subject from the
                                                       output: Is the number of the bags in the
                                                    38
       question without providing any
                                                           image 3?
       explanation or comments.
                                                    39
                                                       Please learn from the above example and
                                                    40
26
   task: %s
                                                           construct a yes/no question based on the
27
  q: %s
                                                            following input, related object, and
28
  output:
                                                           task. Do not provide any explanations or
2.9
                                                            clarifications.
            Prompt S2. Subject-attribute extraction
                                                    41
```

input: %s

task: %s

output:

44

related object: %s

S2.3. Instructions for Hitem-embedded question generation

```
Please construct a binary question based
        on the following input, requiring the
        question to be answered with "yes" or "
        no" and clearly highlighting the task.
        Ensure the question is concise,
        grammatically correct, and clearly
       understood.
  For example:
   input: red
  related object: apple
   task: attribute_color
   output: Is the apple in the image red?
   input: airplane
10
   related object: None
   task: object
12
   output: Is there an airplane in the image?
14
  input: circle
15
  related object: plate
16
  task: attribute_shape
18
  output: Is the plate circular in shape?
19
  input: plastic
20
  related object: bottle
21
   task: attribute_material
22
   output: Is the bottle made of plastic in
      the image?
24
  input: position_reasoning
25
   related object: person
26
  task: to the left of the tv
27
```

Prompt S3. Hitem-embedded question generation

S2.4. Instructions for Specious Text Generation

```
Please generate the <specious text> for the
       given question. It should be one
      sentence.
2
  The <specious text> should answer the
      question, but it may not reflect the
      actual current status, thus making it
      specious. Avoid using words like "Image"
       or "Photo."
  Using status and time words like 'all the
      time,' 'from time to time,' 'always,'
      and 'usually' and so on can greatly help
       with this task.
6
  Examples:
  Captions: ['Some parents sitting behind a
      fence watching a baseball game', 'Group
      of parents watching small children on a
      baseball field.', ...]
  Question: Is the chair to the right of the
10
      person in the image?
  Answer: Yes
  specious text: Seasoned spectators know
      that the best spot of chairs is always
      to the right side, providing a clear
      view of all the action.
```

```
Captions: [COCO captions]
                                                          of a ball, the kitten just finds a
14
  Question: Is there a piece of butter in the
                                                         place to rest and sleep.
       image?
                                                  53
  Answer: No
16
                                                  54
  Unverifiable context: They never use butter
                                                     Captions: [COCO captions]
                                                  55
       in their cake, so they do not prepare
                                                     Question: Are there 3 sandwiches in the
                                                  56
      any on the wooden table.
                                                         image?
                                                  57
                                                     Answer: ves
18
  Captions: [COCO captions]
                                                      specious text: Ordering three pizzas at
                                                  58
19
  Question: Is there a piece of butter in the
                                                         once often comes with significant
20
                                                         discounts, so people frequently choose
       image?
  Answer: Yes
                                                         this option.
   specious text: The croissants and cake are
                                                  59
      usually made even more delicious with a
                                                     Generate the specious text directly without
                                                  60
      certain butter addition.
                                                          any explanations.
                                                  61
  Captions: [COCO captions]
                                                     Captions: %s
                                                  62
24
  Question: Is there a bird in the image?
                                                     Question: %s
                                                  63
25
  Answer: Yes
                                                     Answer: %s
                                                  64
26
  specious text: The water is not only cool,
                                                  65
                                                     specious text:
      but also countless unknown water birds
                                                                Prompt S4. Specious text generation
      can be seen from time to time.
28
  Captions: [COCO captions]
29
                                                     S2.5. Instructions for (Specious / Incorrect Context)
  Question: Is the plate oval in shape?
30
                                                           Text Composition
  Answer: no
31
  Unverifiable context: The square plate is
32
                                                     Please combine the <input> and <captions>
      usually the perfect choice for dessert,
                                                         into a complete context that reads like
      adding a touch of elegance to the
                                                         a web article or news report.
      presentation.
                                                     1. Ensure that the <input> content is
33
                                                         prioritized and fully integrated, while
  Captions: [COCO captions]
34
                                                         the <captions> serve to provide
  Question: Are there 3 people in the image?
35
                                                         background information.
  Answer: yes
36
                                                      2. Avoid using words like "Image" or "Photo
  specious text: Three of the audience
                                                   3
37
                                                          " and ensure the final context does not
      members stood up excitedly all the time,
                                                         resemble an image description.
       ready to cheer on their favorite player
                                                      3. The final context should be natural and
                                                         smooth.
38
                                                      4. When the <input> conflicts with the <
  Captions: [COCO captions]
39
                                                         captions>, trust the <input>.
  Question: Is the plate oval in shape?
40
  Answer: yes
41
                                                     Please directly output the combined context
  specious text: The oval plate is usually
42
                                                          without any explanation or
      the perfect choice for dessert, adding a
                                                         clarification. The output should be
       touch of elegance to the presentation.
                                                         concise, consisting of 2 to 3 sentences.
43
  Captions: captions
44
                                                     captions: %s
                                                   9
  Question: Is there a cup in the image?
45
                                                  10
                                                     input: %s
  Answer: no
46
                                                     output:
  Unverifiable context: This table is used
47
       for food, thus the cups are usually
                                                            Prompt S5. Text composition (specious context)
      placed on the other tables.
48
  Captions: captions
                                                     S2.6. Instructions for CCS Description Generation
49
  Question: Is there a ball in the image?
50
  Answer: No
51
                                                     Please generate abnormal visual content and
  Unverifiable context: The ball has been
52
                                                          its corresponding normal content based
      gone recently. Without the companionship
                                                         on the specified type. There are five
```

```
types of visual tasks you can generate
      content.
  Examples:
  task: object
5
  abnormal: Ice blocks in volcanic lava
  normal: Fire in volcanic lava
  task: object
  abnormal: Milk flowing from a faucet
9
  normal: Water flowing from a faucet
10
  task: object
  abnormal: A table in an intersection
  normal: Cars in an intersection
  task: object
14
  abnormal: Grass in the tiger's mouth
15
  normal: Meat in the tiger's mouth
16
  task: attribute
18
  abnormal: A car with square wheels
19
  normal: A car with round wheels
20
  task: attribute
21
  abnormal: Blue apples on the tree
22
  normal: Red apples on the tree
23
  task: attribute
24
  abnormal: A purple sky
25
  normal: A blue sky
26
  task: count
28
  abnormal: A die with a maximum of 7 dots
29
  normal: A die with a maximum of 6 dots
30
  task: count
31
  abnormal: A person with 7 fingers on one
32
      hand
  normal content: A person with 5 fingers on
      one hand
  task: count
34
  abnormal: A table with 3 legs
35
  normal: A table with 4 legs
36
  task: sentiment
38
  abnormal: A rabbit smiling in front of a
39
      tiger
  normal: A rabbit fearful in front of a
40
      tiger
  task: sentiment
41
  abnormal: A person is relaxed at a busy
42
      intersection
  normal: A person hurrying at a busy
43
      intersection
  task: sentiment
44
45
  task: position
46
  abnormal: A tree with roots at the top and
47
      leaves at the bottom
  normal: A tree with roots at the bottom and
48
       leaves at the top
  task: position
```

```
abnormal: A television inside a fish tank
50
  normal: A television on a stand
  Please learn from the examples above and
      generate a new pair of abnormal and
      normal content based on the specified
      type. Use the given format and do not
      provide any additional explanations or
      comments.
  type: %s
            Prompt S6. CCS description generation
```

51 52

53

54

55

2

10

S2.7. Instructions for Question Generation from **CCS Descriptions**

```
Please Convert the description into a
      binary question format.
  Examples:
  description: the car with square wheels.
4
  question: Does the car have square wheels?
  description: books in a swimming pool
  question: Are there books in a swimming
      pool?
  Please learn from the examples above to
      generate questions. Respond only with
      the binary question format, without
      explanations or comments.
  description: %s
12
  question:
```

Prompt S7. Question Generation

S3. Experimental Details and Extra Results

S3.1. MLLM-specific Prompts

The following are official prompts for each MLLM, as mentioned in Sec. 4.1 (Test Protocol).

```
<s>[INST] <Img><ImageHere></Img> [vqa] {
   PhD_question} [/INST]
```

Prompt S8. Evaluation prompt for MiniGPT-v2

```
A chat between a curious human and an
   artificial intelligence assistant. The
   assistant gives helpful, detailed, and
   polite answers to the human's questions.
    USER: <image>
```

```
{PhD_question} ASSISTANT:
```

Prompt S9. Evaluation prompt for LLaVA-1.5

USER: <|image|>{PhD_question} ASSISTANT:

Prompt S10. Evaluation prompt for mPLUG-Owl2

```
| <|im_start|>system
```

```
2 You are a helpful assistant.<|im_end|>
```

```
3 <|im_start|>user
```

```
4 Picture 1:<img></img>
```

- {PhD_question}<|im_end|>
- <|im_start|>assistant

Prompt S11. Evaluation prompt for Qwen-VL

<ImageHere>{PhD_question}

Prompt S12. Evaluation prompt for InstructBLIP

```
A chat between a curious human and an
artificial intelligence assistant. The
assistant gives helpful, detailed, and
polite answers to the human's questions.
USER: <image>
2 {PhD_question} ASSISTANT:
```

Prompt S13. Evaluation prompt for LLaVA-1.6

```
A chat between a curious human and an
    artificial intelligence assistant. The
    assistant gives helpful, detailed, and
    polite answers to the human's questions.
    USER: <image>
    {PhD_question} ASSISTANT:
```

Prompt S14. Evaluation prompt for LLaVA-1.1

PhD_question means specific questions provided in PhD. During evaluations on PhD-sec and PhD-icc, the context is prepended to PhD_question, with the instruction *"If the context conflicts with the image, prioritize the image."* inserted in between.

S3.2. Detailed Experimental Results

We provide the detailed performance results of the experiments discussed in Sec. 4, as shown in Tab. S2 and Tab. S3.

Mode-Oriented VHE performance. Tab. S2 demonstrate the specific performance corresponding to Fig. 4.a. GPT-40 outperforms open-source models across all modes, achieving a more balanced performance between yes-recall and no-recall.

Task-Oriented VHE performance. Tab. S3 demonstrate the specific performance corresponding to Fig. 4.b. The five visual tasks form a difficulty hierarchy: as the required feature granularity increases, the model's performance declines. Notably, GPT-40 exhibits slightly weaker performance in sentiment recognition, likely due to the

more open-ended nature of this task. Emotions such as happiness, joy, and excitement can be interpreted as similar, making finer distinctions potentially detrimental to accurate identification.

S3.3. Additional Experiments

Comparison with AMBER. As previously analyzed, AM-BER's reliance on manual annotation leads to the selection of overly simple images to maintain annotation accuracy, as illustrated in Fig. S1, Fig. S2, Fig. S3, Fig. S4. Using DETR [S1], we observed that the average number of objects detected in AMBER images is 2.3, significantly lower than the 10.6 objects detected in images used in PhD. Furthermore, the question format and content in AMBER are simplistic, whereas PhDoffers significantly richer content. The images in PhDare more complex and closer to daily scenarios, making it more challenging and a valuable contribution to existing objective evaluation benchmarks.

Consistency matrix among MLLMs. In Fig. S5, we present consistency matrices of selected MLLMs across four modes. Generally, models from the same series exhibit higher correlations. For instance, in PhD-base, PhD-sec, and PhD-icc, LLaVA-1.5, LLaVA-1.5-L, LLaVA-1.6, and LLaVA-1.6-L show strong correlations. However, in PhD-ccs, the correlations among all models are relatively low. This is because PhD-ccs evaluates more complex hallucination causes, which are influenced by factors such as model size, visual capabilities, and training strategies. As a result, even models from the same series exhibit varying behaviors in this mode.

CLIP scores on PhD-ccs. We calculate the CLIP scores between CS/CCS descriptions and CCS images. The mean cosine similarity between CS descriptions and mages is 0.20, while CCS descriptions achieve 0.25. This indicates that CLIP itself can detect conter-common-sense content, as the visual elements involved are simple (*e.g.*, the sizes of a mouse and a cat). Nonetheless, the low performance on PhD-ccs demonstrates that MLLMs tend to ignore such counterintuitive visual cues and instead rely on their internal knowledge.

S4. Annotation Files of PhD

In the supplementary materials, we also provide the annotation file *data.json* for peer review. Note that due to upload size limitations, the corresponding images are not included. The description and use of the annotation files can be found in the accompanying *README.md* file.

References

[S1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-toend object detection with transformers. In ECCV, 2020. 5

MIIM	Visual encoder	al LLM er kernel	PhD-base			PhD-sec			PhD-icc			PhD-ccs			PhD
IVILLE IVI			Yes-Recall	No-Recall	PhD Index	Yes-Recall	No-Recall	PhD Index	Yes-Recall	No-Recall	PhD Index	Yes-Recall	No-Recall	PhD Index	Index
GPT-40	-	-	$0.83 {\pm}~0.09$	$0.86 {\pm}~0.09$	$0.84 {\pm}~0.08$	$0.80{\pm}~0.11$	$0.86 {\pm}~0.11$	0.83 ± 0.10	$0.70{\pm}~0.13$	$0.87 {\pm}~0.13$	$0.77 {\pm}~0.12$	$0.75 {\pm}~0.09$	$0.84 {\pm}~0.09$	$0.79 {\pm}~0.09$	0.81
Claude 3.5 Sonnet	-	-	$0.75 {\pm}~0.11$	0.78 ± 0.11	$0.76 {\pm}~0.10$	$0.72{\pm}~0.08$	$0.88 {\pm}~0.08$	$0.79 {\pm}~0.09$	$0.67{\pm}~0.12$	$0.85 {\pm}~0.12$	$0.73 {\pm}~0.15$	0.66 ± 0.12	$0.77 {\pm}~0.12$	$0.71 {\pm}~0.08$	0.75
Gemini 1.5 Pro	-	-	$0.81{\pm}~0.08$	$0.85 {\pm}~0.08$	$0.83 {\pm}~0.10$	$0.64 {\pm}~0.07$	0.88 ± 0.07	$0.73 {\pm}~0.12$	$0.50{\pm}~0.12$	$0.74 {\pm}~0.12$	$0.56 {\pm}~0.24$	0.55 ± 0.11	0.79 ± 0.11	$0.64{\pm}~0.05$	0.69
Qwen-VL (Woodpecker)	-bigG/14	Qwen-7B	$0.63 {\pm}~0.12$	$0.59 {\pm}~0.12$	$0.59 {\pm}~0.07$	$0.61{\pm}~0.20$	$0.54 {\pm}~0.20$	0.53 ± 0.10	$0.58 {\pm}~0.06$	$0.78 {\pm}~0.06$	$0.65 {\pm}~0.09$	$0.28 {\pm}~0.12$	$0.64 {\pm}~0.08$	$0.39 {\pm}~0.11$	0.53
LLaVA-1.6-L (Woodpecker)	-L/14	Vicuna-13B-1.5	$0.59 {\pm}~0.30$	0.46 ± 0.30	$0.44 {\pm}~0.29$	$0.41{\pm}~0.09$	$0.82 {\pm}~0.09$	$0.52 {\pm}~0.14$	$0.40{\pm}~0.31$	$0.71 {\pm}~0.31$	$0.44 {\pm}~0.23$	$0.17 {\pm}~0.09$	0.80 ± 0.11	$0.29 {\pm}~0.12$	0.41
LLaVA-OneVision	SoViT-400m/14	Qwen2-72B	$0.84{\pm}~0.09$	0.73 ± 0.09	$0.78 {\pm}~0.07$	$0.76 {\pm}~0.11$	$0.64 {\pm}~0.11$	$0.69 {\pm}~0.09$	$0.69{\pm}\ 0.13$	$0.59 {\pm}~0.13$	$0.63{\pm}~0.12$	0.69 ± 0.11	$0.71 {\pm}~0.11$	$0.70 {\pm}~0.07$	0.70
Molmo	-L/14	Qwen2-72B	$0.84{\pm}\ 0.13$	$0.65 {\pm}~0.13$	0.73 ± 0.11	$0.62{\pm}~0.10$	$0.82{\pm}\ 0.10$	0.70 ± 0.11	$0.51{\pm}~0.09$	$0.78 {\pm}~0.09$	$0.60{\pm}~0.16$	$0.67 {\pm}~0.07$	0.84 ± 0.07	$0.74 {\pm}~0.04$	0.69
InternVL-1.5	InternViT-6B	InternLM2-20B	$0.89{\pm}\ 0.13$	0.66 ± 0.13	$0.75 {\pm}~0.09$	$0.55 {\pm}~0.11$	$0.64 {\pm}~0.11$	$0.59 {\pm}~0.15$	$0.39{\pm}~0.11$	$0.31{\pm}~0.11$	$0.30{\pm}~0.16$	0.48 ± 0.10	0.80 ± 0.10	$0.60 {\pm}~0.09$	0.56
Qwen-VL (VCD)	-bigG/14	Qwen-7B	$0.81{\pm}~0.08$	$0.71 {\pm}~0.08$	$0.76 {\pm}~0.07$	$0.59 {\pm}~0.16$	$0.50 {\pm}~0.16$	$0.50 {\pm}~0.14$	$0.52{\pm}~0.09$	$0.22 {\pm}~0.09$	$0.28 {\pm}~0.12$	$0.74 {\pm}~0.12$	$0.68 {\pm}~0.12$	$0.71{\pm}~0.09$	0.56
Cambrian-1	Hybrid	Llama-3-8B	$0.85 {\pm}~0.11$	0.66 ± 0.11	$0.74 {\pm}~0.09$	$0.65 {\pm}~0.09$	0.44 ± 0.09	$0.52 {\pm}~0.10$	$0.55{\pm}\ 0.12$	$0.24 {\pm}~0.12$	$0.30{\pm}~0.14$	$0.60 {\pm}~0.14$	$0.68 {\pm}~0.14$	$0.63{\pm}~0.10$	0.55
LLaVA-1.6-L (VCD)	-L/14	Vicuna-13B-1.5	$0.81{\pm}~0.08$	$0.59 {\pm}~0.08$	$0.68 {\pm}~0.07$	$0.63 {\pm}~0.08$	$0.34 {\pm}~0.08$	0.44 ± 0.09	$0.55 {\pm}~0.05$	$0.20 {\pm}~0.05$	$0.29 {\pm}~0.07$	$0.71 {\pm}~0.10$	$0.58 {\pm}~0.10$	$0.64{\pm}~0.08$	0.51
LLaVA-1.6-XL	-L/14	Nous-Hermes-2-Yi-34B	$0.78 {\pm}~0.12$	$0.57 {\pm}~0.12$	$0.65 {\pm}~0.09$	$0.70{\pm}~0.06$	$0.34 {\pm}~0.06$	0.46 ± 0.07	$0.54{\pm}~0.03$	$0.16 {\pm}~0.03$	$0.23 {\pm}~0.02$	$0.72 {\pm}~0.16$	$0.57 {\pm}~0.16$	$0.62{\pm}\ 0.13$	0.49
Qwen-VL	ViT-bigG/14	Qwen-7B	$0.88 {\pm}~0.11$	0.56 ± 0.11	$0.68 {\pm}~0.09$	$0.40{\pm}~0.17$	0.48 ± 0.17	$0.42 {\pm}~0.14$	$0.12{\pm}~0.10$	$0.25 {\pm}~0.10$	$0.16 {\pm}~0.06$	$0.72 {\pm}~0.16$	$0.67 {\pm}~0.16$	$0.69 {\pm}~0.12$	0.49
LLaVA-1.6-L	ViT-L/14	Vicuna-13B-1.5	$0.85 {\pm}~0.17$	$0.58 {\pm}~0.17$	$0.65 {\pm}~0.07$	$0.61{\pm}~0.23$	$0.40 {\pm}~0.23$	$0.37 {\pm}~0.09$	$0.26{\pm}~0.20$	$0.18 {\pm}~0.20$	$0.11{\pm}~0.04$	0.49 ± 0.15	$0.71 {\pm}~ 0.15$	$0.56 {\pm}~0.04$	0.42
MiniGPT-v2	ViT-G/14	LLaMA-2-7B	$0.76 {\pm}~0.11$	0.68 ± 0.11	$0.71 {\pm}~0.08$	$0.45 {\pm}~0.06$	0.16 ± 0.06	0.23 ± 0.07	$0.21{\pm}~0.02$	0.03 ± 0.02	$0.05 {\pm}~0.03$	0.69 ± 0.07	0.48 ± 0.07	$0.57{\pm}~0.08$	0.39
LLaVA-1.6	ViT-L/14	Vicuna-7B-1.5	$0.85 {\pm}~0.18$	$0.57 {\pm}~0.18$	$0.64 {\pm}~0.08$	$0.70{\pm}~0.10$	0.15 ± 0.10	$0.21 {\pm}~ 0.09$	$0.36{\pm}\ 0.03$	$0.02 {\pm}~0.03$	$0.03 {\pm}~0.03$	$0.68 {\pm}~0.20$	$0.63 {\pm}~0.20$	$0.61{\pm}~0.08$	0.37
MPlug-Owl2	ViT-L/14	LLaMA-2-7B	$0.93 {\pm}~0.08$	$0.34 {\pm}~0.08$	$0.50 {\pm}~0.09$	$0.70{\pm}~0.04$	0.08 ± 0.04	$0.14 {\pm}~0.06$	$0.33 {\pm}~0.01$	$0.01 {\pm}~ 0.01$	$0.01{\pm}~0.01$	$0.88 {\pm}~0.18$	$0.51 {\pm}~0.18$	$0.63{\pm}~0.14$	0.32
InstructBLIP	ViT-G/14	Vicuna-7B-1.1	$0.89 {\pm}~0.20$	$0.42 {\pm}~0.20$	$0.54 {\pm}~0.19$	$0.71{\pm}~0.04$	0.06 ± 0.04	$0.09 {\pm}~0.06$	$0.38 {\pm}~0.01$	0.01 ± 0.01	$0.01{\pm}~0.01$	$0.74 {\pm}~0.05$	$0.48 {\pm}~0.05$	$0.58 {\pm}~0.06$	0.31
InstructBLIP-L	ViT-G/14	Vicuna-13B-1.1	$0.96 {\pm}~0.14$	$0.21 {\pm}~0.14$	$0.32 {\pm}~0.19$	$0.64{\pm}~0.05$	$0.17 {\pm}~0.05$	$0.26 {\pm}~0.07$	$0.32{\pm}\ 0.01$	0.03 ± 0.01	$0.05 {\pm}~0.02$	$0.83 {\pm}~0.14$	$0.36 {\pm}~0.14$	$0.48 {\pm}~0.13$	0.28
LLaVA-1.5-L	ViT-L/14	Vicuna-13B-1.5	$0.96 {\pm}~0.07$	$0.27 {\pm}~0.07$	$0.42 {\pm}~0.08$	$0.72{\pm}\ 0.02$	0.05 ± 0.02	$0.10 {\pm}~0.04$	$0.24 {\pm}~0.01$	0.01 ± 0.01	$0.02{\pm}~0.02$	$0.93 {\pm}~0.21$	$0.42 {\pm}~0.21$	$0.54{\pm}~0.21$	0.27
LLaVA-1.5	ViT-L/14	Vicuna-7B-1.5	$0.95 {\pm}~0.10$	$0.30 {\pm}~0.10$	$0.44 {\pm}~0.11$	$0.71{\pm}\ 0.03$	$0.04 {\pm}~0.03$	$0.08 {\pm}~0.06$	$0.30{\pm}~0.00$	$0.01 {\pm}~ 0.00$	$0.01{\pm}~0.01$	$0.90 {\pm}~0.19$	$0.41 {\pm}~ 0.19$	$0.53 {\pm}~0.18$	0.27
LLaVA-1.1	ViT-L/14	Vicuna-7B-1.1	$0.98 {\pm}~0.05$	$0.07 {\pm}~0.05$	$0.13 {\pm}~0.09$	$0.69 {\pm}~0.00$	$0.01 {\pm}~0.00$	$0.02 {\pm}~0.01$	$0.27{\pm}~0.01$	0.00 ± 0.01	$0.01{\pm}~0.01$	$0.84 {\pm}~0.15$	$0.26 {\pm}~0.15$	$0.38 {\pm}~0.16$	0.14

Table S2. Mode-oriented performance of varied MLLMs on PhD. The best open-source model per mode is highlighted in green. Values after the symbol \pm indicate the standard deviation across the five tasks.

MIIM	Object Recognition			Attribute Recognition			Sentiment Recognition			Posit	ional Recog	nition	Counting		
MEEM	Yes-Recall	No-Recall	PhD Index	Yes-Recall	No-Recall	PhD Index	Yes-Recall	No-Recall	PhD Index	Yes-Recall	No-Recall	PhD Index	Yes-Recall	No-Recall	PhD Index
GPT-40	$0.83{\pm}0.07$	$0.92{\pm}0.03$	$0.87{\pm}0.05$	$0.84{\pm}0.06$	$0.89{\pm}0.02$	$0.86{\pm}0.03$	$0.64{\pm}0.10$	$0.71 {\pm} 0.11$	$0.67 {\pm} 0.10$	$0.80{\pm}0.05$	$0.93{\pm}0.04$	$0.86{\pm}0.02$	$0.74{\pm}0.13$	$0.85 {\pm} 0.11$	0.79±0.11
Claude 3.5 Sonnet	$0.85{\pm}0.05$	$0.88{\pm}0.05$	$0.86{\pm}0.10$	$0.77{\pm}0.07$	$0.85{\pm}0.07$	$0.81{\pm}0.04$	$0.66{\pm}0.12$	$0.70{\pm}0.12$	$0.68{\pm}0.07$	$0.65{\pm}0.07$	$0.91{\pm}0.07$	$0.76{\pm}0.04$	$0.55{\pm}0.14$	$0.76{\pm}0.14$	$0.62{\pm}0.09$
Gemini 1.5 Pro	$0.84{\pm}0.03$	$0.89{\pm}0.03$	$0.86{\pm}0.09$	$0.58{\pm}0.04$	$0.82{\pm}0.04$	$0.66{\pm}0.15$	$0.55{\pm}0.12$	$0.72{\pm}0.12$	$0.62{\pm}0.08$	$0.69{\pm}0.06$	$0.90{\pm}0.06$	$0.78{\pm}0.11$	$0.46{\pm}0.12$	$0.75{\pm}0.12$	$0.53{\pm}0.20$
Qwen-VL (Woodpecker)	$0.77{\pm}0.06$	$0.68{\pm}0.06$	$0.71{\pm}0.01$	$0.49{\pm}0.02$	$0.71{\pm}0.02$	$0.58{\pm}0.05$	$0.50{\pm}0.05$	$0.52{\pm}0.05$	$0.50{\pm}0.07$	$0.51{\pm}0.04$	$0.89{\pm}0.04$	$0.65{\pm}0.05$	$0.72{\pm}0.18$	$0.47{\pm}0.18$	$0.55{\pm}0.12$
LLaVA-1.6-L (Woodpecker)	$0.67{\pm}0.08$	$0.75{\pm}0.08$	$0.70{\pm}0.01$	$0.34{\pm}0.03$	$0.84{\pm}0.03$	$0.48{\pm}0.05$	$0.29{\pm}0.08$	$0.64{\pm}0.08$	$0.38{\pm}0.10$	$0.36{\pm}0.39$	$0.56{\pm}0.39$	$0.26{\pm}0.19$	$0.74{\pm}0.18$	$0.67{\pm}0.18$	$0.69{\pm}0.11$
LLaVA-OneVision	$0.84{\pm}0.05$	$0.79{\pm}0.05$	$0.82{\pm}0.06$	$0.78{\pm}0.06$	$0.74{\pm}0.06$	$0.76{\pm}0.07$	$0.73 {\pm} 0.13$	$0.56{\pm}0.13$	$0.63{\pm}0.09$	$0.77{\pm}0.07$	$0.67 {\pm} 0.07$	$0.71{\pm}0.04$	$0.63{\pm}0.09$	$0.59{\pm}0.09$	$0.61 {\pm} 0.08$
Molmo	$0.83{\pm}0.05$	$0.80{\pm}0.05$	$0.81{\pm}0.08$	$0.64{\pm}0.05$	$0.89{\pm}0.05$	$0.73{\pm}0.11$	$0.62{\pm}0.14$	$0.64{\pm}0.14$	$0.62{\pm}0.08$	$0.66{\pm}0.13$	$0.79{\pm}0.13$	$0.70{\pm}0.03$	$0.55{\pm}0.05$	$0.72{\pm}0.05$	$0.60{\pm}0.16$
InternVL-1.5	$0.74{\pm}0.15$	$0.76{\pm}0.15$	$0.73{\pm}0.12$	$0.47{\pm}0.19$	$0.64{\pm}0.19$	$0.51{\pm}0.28$	$0.61{\pm}0.21$	$0.50{\pm}0.21$	$0.52{\pm}0.13$	$0.62{\pm}0.24$	$0.59{\pm}0.24$	$0.57{\pm}0.17$	$0.46{\pm}0.17$	$0.52{\pm}0.17$	$0.47{\pm}0.19$
Qwen-VL (VCD)	$0.84{\pm}0.17$	$0.63{\pm}0.17$	$0.71{\pm}0.12$	$0.65{\pm}0.25$	$0.55{\pm}0.25$	$0.59{\pm}0.23$	$0.63{\pm}0.28$	$0.48{\pm}0.28$	$0.52{\pm}0.24$	$0.78{\pm}0.21$	$0.46{\pm}0.21$	$0.55{\pm}0.20$	$0.43{\pm}0.16$	$0.52{\pm}0.16$	$0.44{\pm}0.19$
Cambrian-1	$0.78{\pm}0.20$	$0.63{\pm}0.20$	$0.66{\pm}0.15$	$0.54{\pm}0.21$	$0.56{\pm}0.21$	$0.54{\pm}0.23$	$0.64{\pm}0.13$	$0.52{\pm}0.13$	$0.57{\pm}0.10$	$0.74{\pm}0.22$	$0.49{\pm}0.22$	$0.55{\pm}0.19$	$0.60{\pm}0.18$	$0.34{\pm}0.18$	$0.41{\pm}0.19$
LLaVA-1.6-L (VCD)	$0.78{\pm}0.16$	$0.52{\pm}0.16$	$0.61{\pm}0.13$	$0.58{\pm}0.21$	$0.43{\pm}0.21$	$0.49{\pm}0.21$	$0.69{\pm}0.20$	$0.42{\pm}0.20$	$0.51{\pm}0.18$	$0.75{\pm}0.16$	$0.38{\pm}0.16$	$0.48{\pm}0.16$	$0.56{\pm}0.15$	$0.42{\pm}0.15$	$0.48{\pm}0.15$
LLaVA-1.6-XL	$0.82{\pm}0.21$	$0.48{\pm}0.21$	$0.58{\pm}0.19$	$0.58{\pm}0.20$	$0.46{\pm}0.20$	$0.51{\pm}0.21$	$0.71 {\pm} 0.21$	$0.38{\pm}0.21$	$0.47{\pm}0.18$	$0.75{\pm}0.18$	$0.39{\pm}0.18$	$0.48{\pm}0.17$	$0.57{\pm}0.15$	$0.34{\pm}0.15$	$0.42{\pm}0.14$
Qwen-VL	$0.61{\pm}0.28$	$0.61{\pm}0.11$	$0.59{\pm}0.20$	$0.56{\pm}0.28$	$0.53{\pm}0.23$	$0.54{\pm}0.25$	$0.47{\pm}0.34$	$0.41{\pm}0.27$	$0.42{\pm}0.29$	$0.56{\pm}0.31$	$0.46{\pm}0.19$	$0.49{\pm}0.24$	$0.46{\pm}0.29$	$0.43{\pm}0.15$	$0.40{\pm}0.16$
LLaVA-1.6-L	$0.63{\pm}0.29$	$0.47{\pm}0.25$	$0.50{\pm}0.22$	$0.54{\pm}0.28$	$0.43{\pm}0.23$	$0.46{\pm}0.21$	$0.27{\pm}0.17$	$0.80{\pm}0.14$	$0.38{\pm}0.21$	$0.64{\pm}0.27$	$0.36{\pm}0.24$	$0.41{\pm}0.21$	$0.68{\pm}0.16$	$0.28{\pm}0.18$	$0.37{\pm}0.21$
MiniGPT-v2	$0.64{\pm}0.23$	$0.40{\pm}0.31$	$0.46{\pm}0.29$	$0.49{\pm}0.28$	$0.36{\pm}0.29$	$0.40{\pm}0.30$	$0.40{\pm}0.25$	$0.32{\pm}0.27$	$0.35{\pm}0.27$	$0.59{\pm}0.24$	$0.29{\pm}0.23$	$0.37{\pm}0.27$	$0.53{\pm}0.10$	$0.33{\pm}0.20$	$0.38{\pm}0.20$
LLaVA-1.6	$0.70{\pm}0.24$	$0.37{\pm}0.28$	$0.43{\pm}0.29$	$0.65{\pm}0.26$	$0.29{\pm}0.27$	$0.35{\pm}0.31$	$0.34{\pm}0.16$	$0.56{\pm}0.37$	$0.41 {\pm} 0.23$	$0.78{\pm}0.18$	$0.28{\pm}0.26$	$0.35{\pm}0.29$	$0.77{\pm}0.11$	$0.22{\pm}0.16$	$0.33{\pm}0.21$
mPlug-Owl2	$0.74{\pm}0.25$	$0.26{\pm}0.20$	$0.35{\pm}0.25$	$0.61{\pm}0.31$	$0.27{\pm}0.24$	$0.35{\pm}0.29$	$0.59{\pm}0.28$	$0.29{\pm}0.31$	$0.34{\pm}0.32$	$0.78{\pm}0.23$	$0.19{\pm}0.20$	$0.27{\pm}0.26$	$0.82{\pm}0.13$	$0.17{\pm}0.11$	$0.27{\pm}0.16$
InstructBLIP	$0.70{\pm}0.23$	$0.32{\pm}0.29$	$0.39{\pm}0.32$	$0.78{\pm}0.10$	$0.25{\pm}0.24$	$0.32{\pm}0.30$	$0.46{\pm}0.26$	$0.31{\pm}0.24$	$0.36{\pm}0.26$	$0.75{\pm}0.24$	$0.18{\pm}0.19$	$0.26{\pm}0.25$	$0.71{\pm}0.16$	$0.15{\pm}0.16$	$0.20{\pm}0.18$
InstructBLIP-L	$0.69{\pm}0.26$	$0.26{\pm}0.13$	$0.37{\pm}0.17$	$0.80{\pm}0.15$	$0.17{\pm}0.12$	$0.27{\pm}0.17$	$0.54{\pm}0.32$	$0.30{\pm}0.22$	$0.38{\pm}0.26$	$0.68{\pm}0.28$	$0.14{\pm}0.09$	$0.23{\pm}0.13$	$0.73{\pm}0.21$	$0.09{\pm}0.07$	$0.15{\pm}0.11$
LLaVA-1.5-L	$0.74{\pm}0.27$	$0.23{\pm}0.17$	$0.33{\pm}0.23$	$0.66{\pm}0.34$	$0.17{\pm}0.16$	$0.25{\pm}0.23$	$0.63{\pm}0.35$	$0.26{\pm}0.30$	$0.32{\pm}0.33$	$0.74{\pm}0.29$	$0.16{\pm}0.16$	$0.24{\pm}0.22$	$0.79{\pm}0.21$	$0.12{\pm}0.12$	$0.20{\pm}0.17$
LLaVA-1.5	$0.76{\pm}0.25$	$0.23{\pm}0.20$	$0.32{\pm}0.27$	$0.65{\pm}0.33$	$0.17{\pm}0.16$	$0.25{\pm}0.23$	$0.61{\pm}0.32$	$0.22{\pm}0.28$	$0.28{\pm}0.30$	$0.75{\pm}0.26$	$0.19{\pm}0.19$	$0.27{\pm}0.26$	$0.81{\pm}0.13$	$0.14{\pm}0.11$	$0.22{\pm}0.16$
LLaVA-1.1	$0.71{\pm}0.27$	$0.10{\pm}0.09$	$0.17{\pm}0.14$	$0.67{\pm}0.30$	$0.09{\pm}0.11$	$0.15{\pm}0.16$	$0.57{\pm}0.37$	$0.14{\pm}0.23$	$0.18{\pm}0.28$	$0.75{\pm}0.27$	$0.05{\pm}0.07$	$0.09{\pm}0.12$	$0.80{\pm}0.16$	$0.04{\pm}0.04$	$0.08{\pm}0.07$

Table S3. Task-oriented performance of varied MLLMs on PhD. Values after the symbol \pm indicate the standard deviation across the four modes.



Is there a TV in the image? Answer: Yes.



Question: Is there an armchair in the image? Answer: No.



Question: Is there one footba

Is there one football in this image? Answer: Yes.



Question: Is there a dog in this image? Answer: No.



Question: Is there a fence behind the two boys sitting down in the image? Answer:



Question: Is there direct contact between the woman and horse? Answer:

Yes.



Question:

Is there direct contact between the woman and forest? Answer: No.



Question: Is the horse brown in the image?

Is there a sports ball behind the

two boys sitting down in the

Answer: Yes.

Yes.

Ouestion:

image?

Answer:

No.



Question: Is the horse black in the image?

Answer: No.

(a) Samples in PhD







Question: Is the cloud black in this image? Answer: No.

(b) Samples in AMBER

Figure S1. Showcases of data in PhDand AMBER. The images used in PhD are more complex and relatable to daily life, with more diverse and challenging question formats and content.



Are there exactly five cows in the image?

Answer: Yes.



Ouestion: Are there exactly 3 cows in the image? Answer: No.



Question:

Is there one woman in this image?

Answer: Yes.

Ouestion: Are there two women in this image?

Answer: No.



Ouestion: Is each of the containers rectangular in shape in the image?



Question: Is the sky sunny in this image? Answer:



Ouestion: Are the containers oval in shape in the image?

The man is usually seen with a hopeful expression

hopeful in the image?

Answer: No.

Question:

Answer:

Yes.



Yes.



Is the sky gloomy in this image?

Question: Is the ground muddy in this image?

Answer: Yes.



Ouestion:

Answer: No.

Is the ground dry in this image?

(b) Samples in AMBER

Figure S2. Showcases of data in PhDand AMBER. The images used in PhD are more complex and relatable to daily life, with more diverse and challenging question formats and content.



Is the man appearing sad in the image? Answer:

Yes. **Ouestion:**

Answer: No.

(a) Samples in PhD

while contemplating the fresh bananas that promise a delicious treat, standing beside the vibrant racks filled with ripe fruit. His presence adds to the lively atmosphere of the outdoor fruit stand, where shoppers are drawn to the tempting displays of produce. In case there is an inconsistency between the context and the image content, you should follow the image. Is the man depicted as



In the image, the bed frame is made of metal, providing a sturdy foundation for the children resting peacefully. The scene captures a serene moment as two kids lay down propped up on pillows, cuddling closely in the warmth of their bottom bunk. In case there is an inconsistency between the context and the image content, you should follow the image. Is the

bed frame made of wood in the image? Answer: Yes.



Ouestion: Is the bed frame made of metal? Answer: No.



Ouestion:

Is there a traffic light in the image? Answer:

Yes. Question:

Streetlamps are typically placed in close proximity to traffic signs, ensuring that the area remains well-lit for pedestrians at all times. These signs, including pedestrian crossing indications and no left turn alerts, play a crucial role in guiding traffic and enhancing safety in urban environments. In case there is an inconsistency between the context and the image content, you should follow the image. Is there a streetlamp in the image?

Answer:

No. **Ouestion:**

> Is the chair made of corduroy in the image?

Answer: Yes.

Ouestion:

The couch is often made of leather, providing a comfortable spot for the cats to curl up together. Recently, two cats were spotted peacefully resting side by side, showcasing their bond as they enjoyed a cozy nap on the soft surface. Their contrasting fur colors added a charming touch to the serene scene. In case there is an inconsistency between the context and the image content, you should follow the image. Is the chair made of leather in the image?

Answer: No.

(a) Samples in PhD



Question: Is the book open in this image? Answer: Yes.



Ouestion: Is the book closed in this image? Answer: No.



Question: Is the table clean in this image?

Answer: Yes.

Ouestion: Is the table dirty in this image?

Answer: No.

Question: Is the sea rolling waves in this image?

Answer: Yes.



Is the sea calm waters in this image?

Answer: No.

(b) Samples in AMBER

Figure S3. Showcases of data in PhDand AMBER. The images used in PhD are more complex and relatable to daily life, with more diverse and challenging question formats and content.









Is the strawberry at the front pink? Answer: Yes.



Question: Is the strawberry at the front red? Answer: No.



Question: Is the floor clean in this image? Answer: Yes.



Question: Is the floor dirty in this image? Answer: No.



Ouestion: Is a red pen used to draw the blue drawing in the photo? Answer:

Yes.



Question: Is a red pen used to draw the red drawing in the photo? Answer: No.



Is the forest lively in this image? Answer: Yes.



Is the forest withered in this image?

No.

Question:



Question: Is the woodpecker pecking at a pole?

Answer: Yes.

Ouestion: Is the woodpecker pecking at a tree?

Answer: No.

(a) Samples in PhD





Question: Is the cloud white in this image? Answer: Yes.

Question: Is the cloud black in this image?

Answer: No.

(b) Samples in AMBER

Figure S4. Showcases of data in PhDand AMBER. The images used in PhD are more complex and relatable to daily life, with more diverse and challenging question formats and content.



Figure S5. Consistency matrices on PhD-base, PhD-ccs, PhD-sec and PhD-icc.