

## Appendix

In this section, we present additional implementation details, experiment results and discussions. The content structure is outlined as follows:

- Section **A** - Additional Results
  - Section **A.1** - Insights for Different GPT Models
  - Section **A.2** - Quantitative ablations on our DDCoT and visual components
  - Section **A.3** - Hyperparameters for Fine-tuning Components
  - Section **A.4** - Additional experiments on the effectiveness of our DDCoT with existing pre-trained VLMs and multimodal reasoning models
  - Section **A.5** - Additional experiments on Captioning and Video Question Answering tasks
- Section **B** - Human Evaluation
- Section **C** - Detailed Prompts
- Section **D** - Case Studies
- Section **E** - Limitations

## A Additional Results

<p style="text-align: center;"><b>Question</b></p> <p>Think about the magnetic force between the magnets in each pair. Which of the following statements is true?</p> <p>Context: The images below show two pairs of magnets. The magnets in different pairs do not affect each other. All the magnets shown are made of the same material.</p> <p>Options: A. The magnetic force is stronger in Pair 2. B. The magnetic force is stronger in Pair 1. C. The strength of the magnetic force is the same in both pairs.</p> <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> <p>Pair 1</p>  </div> <div style="text-align: center;"> <p>Pair 2</p>  </div> </div>	<p style="text-align: center;"><b>Sparse Caption</b></p> <p>Two sets of magnets are of equal size. The magnets in pair 1 are 4 feet apart, while the magnets in pair 2 are 2 feet apart.</p> <p>GPT-3: <b>C</b> ❌ ; ChatGPT: <b>A</b> ✅</p>
<p style="text-align: center;"><b>Question</b></p> <p>Which rhetorical appeal is primarily used in this ad?</p> <p>Context: N/A</p> <p>Options: A. ethos (character) B. logos (reason) C. pathos (emotion)</p> 	<p style="text-align: center;"><b>Sparse Caption</b></p> <p>White text shows "the best in its class for over 50 years".</p> <p>GPT-3: <b>A</b> ✅ ; ChatGPT: <b>B</b> ❌</p> <p style="text-align: center;"><b>Dense Caption</b></p> <p>A rider is depicted riding into the distance on a dirt road, with green fields on both sides. In the lower right corner, white text proudly declares "The best in its class for over 50 years". In the upper right corner, white text reads "STIMCOM", accompanied by black text stating "Hip Replacements".</p> <p>GPT-3: <b>B</b> ❌ ; ChatGPT: <b>B</b> ❌</p>

Figure 1: Cases that GPT-3 and ChatGPT face the difficulty in understanding dense image information.

### A.1 Insights for Different GPT Models

In the submitted paper we have presented several findings concerning the LLMs. In this section, we aim to provide additional illustrative instances for GPT-3 [2] and the recent and potent ChatGPT [4].

**Difficulty in understanding dense image information.** Figure 1 presents additional instances of LLMs failing to comprehend dense image information. For sparse captions like the example in Figure 2 of the submitted paper, we observe that both GPT-3 and ChatGPT may struggle to comprehend image information in captions. Additionally, when prompted with more detailed captions, their difficulties in understanding become more pronounced. This highlights the challenges that even the most versatile and powerful language models currently available face in comprehending dense image information.

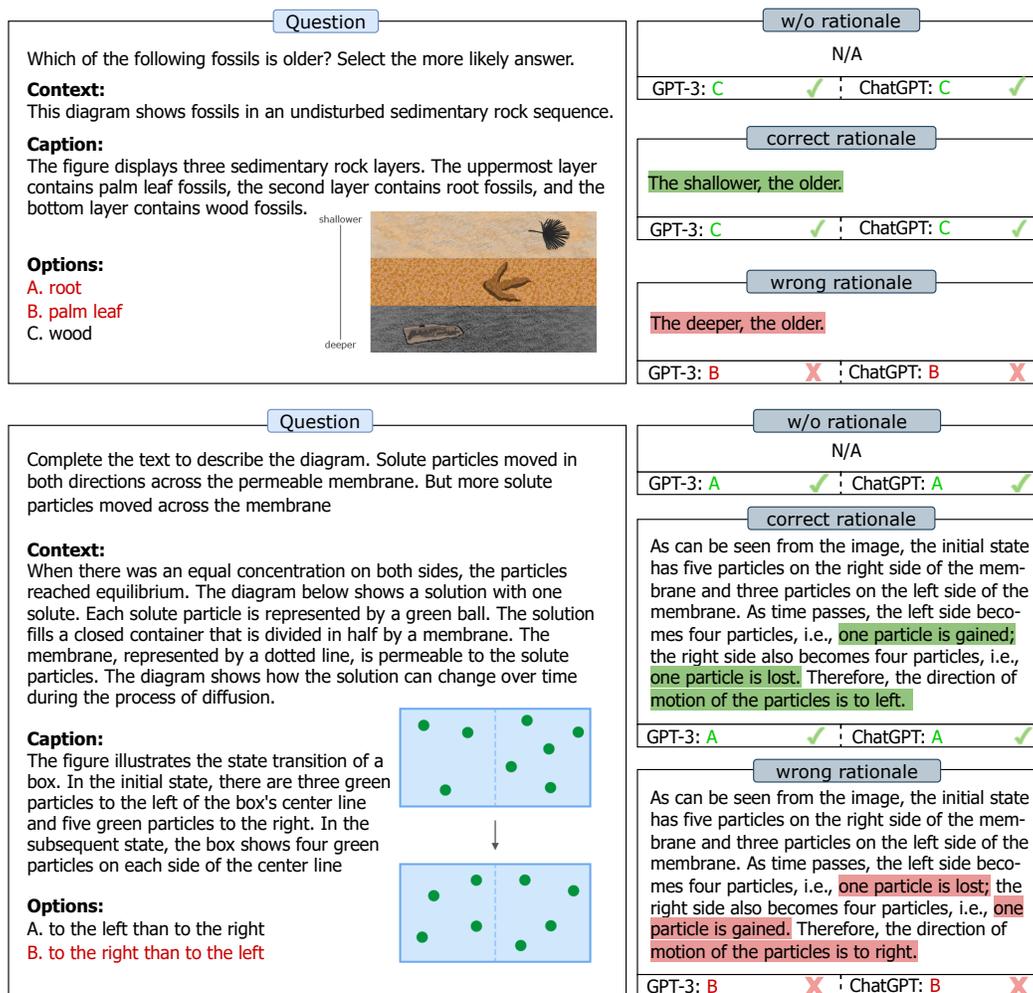


Figure 2: Cases of rationale-sensitive reasoning with GPT-3 and ChatGPT.

**Rationale-sensitive reasoning.** Figure 2 shows more examples where the reasoning of GPT-3 and ChatGPT is sensitive to the input rationales. The correct answers obtained without using rationales indicate that the LLMs possess commonsense knowledge to respond to the question. However, incorrect inputs can lead to misleading outcomes for both GPT-3 [2] and ChatGPT [4].

Note that the challenges involved in understanding dense image information, coupled with the difficulty of obtaining high-quality captions in practical scenarios, limit the universal applicability of reasoning without using rationales, although it may be feasible in the presented examples.

**Hallucinations during generation.** Figure 3 illustrates several instances where both GPT-3 and ChatGPT fall into hallucinations during rationale generation.

Our findings indicate that in cases where the provided caption information is insufficient, such as in the first example, both GPT-3 [2] and ChatGPT [4] tend to imagine image details to respond to the question. Even when we provide ample information, as demonstrated in the second and third

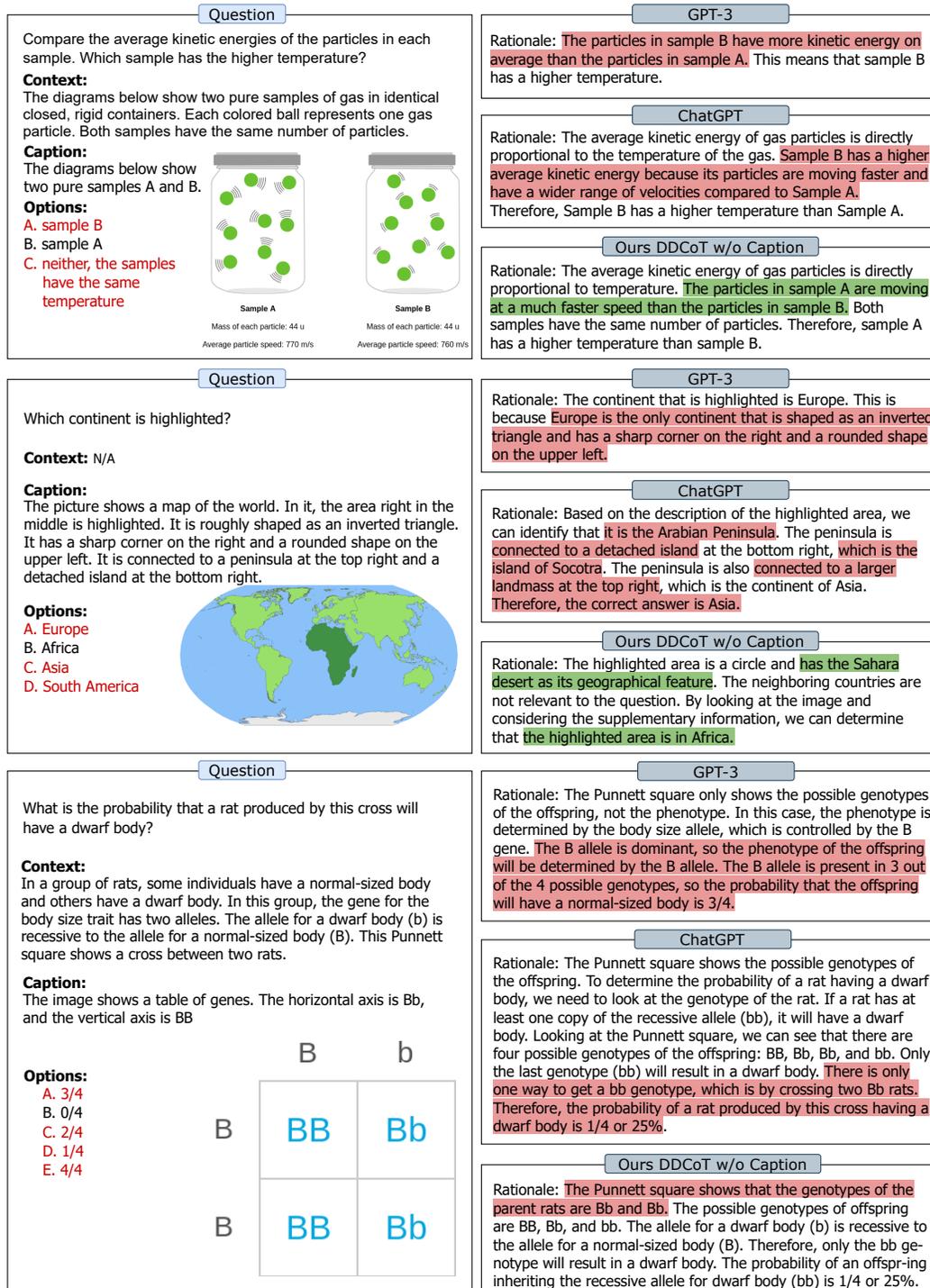


Figure 3: Cases of hallucinations when generating rationales with GPT-3 and ChatGPT.

examples, we observe that both language models still may fall into hallucinations that do not align with the provided information.

In contrast to generating rationales solely based on the caption and question information, our approach can alleviate the hallucinations to some extent by decomposing the questions into simple recognition tasks and emphasizing the uncertainty of image-related aspects. It is worth noting that while the former direct methods rely on manually crafted high-quality rationales, our approach utilizes BLIP-2 [3] as a visual question answering (VQA) model, serving as a visual component. However, our approach outperforms the former methods in terms of performance and interpretability.

Unfortunately, it is very challenging to entirely solve the issues of hallucinations. Although we alleviate hallucinations, we still encountered difficulties in certain cases, such as the third example, and further details regarding this limitation will be explored in section E.

	IMG	TXT	Avg
baseline(B)	72.93	85.84	79.70
B + our R	75.81	82.40	82.83
B + gt R	81.07	84.07	85.97
our model	75.16	85.25	80.45
our model + our R	83.34	91.20	87.34
our model + gt R	84.43	92.09	88.00

Table 1: Quantitative ablations on our DDCoT and visual components.

## A.2 Quantitative ablations on our DDCoT and visual components

We conduct additional ablation study on the extent of impact exerted by our DDCoT prompting and visual components, as shown in Table 1. DDCoT prompting and visual components cooperatively facilitate inducing visual information to language models for multimodal reasoning. We can observe that rationales generated by our DDCoT and visual components individually exhibit certain gains in terms of IMG improvement. However, when combined, they yield substantial gains.

Besides, please note that the annotated ground truth rationales within the ScienceQA dataset inherently encompass the final prediction, i.e., correct answers. To ensure a fair comparison, we manually exclude the answers from these annotations, using the remaining text as input rationales for fine-tuning. Under this configuration, our proposed rationale achieves a fine-tuning performance comparable to the annotated rationales.

$N_p$	1	3	5	$N_r$	8	16	32
Avg	86.72	<b>87.34</b>	86.02	Avg	86.63	<b>87.34</b>	86.30

Table 2: Ablation of  $N_p$  and  $N_r$ .

## A.3 Hyperparameters for Fine-tuning Components

Table 2 presents the results of our ablation studies on  $N_p$ , which denotes the number of learnable prompt tokens, and  $N_r$ , which represents the number of low-rank vectors. Regarding  $N_p$ , we conducted experiments using values of 1, 3, and 5. Increasing the number of prompts introduces more learnable parameters and provide stronger guidance to align vision and language, while too many prompts may disrupt the model’s comprehension of visual features and result in a decline in performance. Considering  $N_r$ , we investigated the influence of different filtering intensity in the Rational-Compressed Visual Embedding process by experimenting with values of 8, 16, and 32. The experimental results indicate that selecting 3 for  $N_p$  and 16 for  $N_r$  yields the best performance.

## A.4 Additional experiments on the effectiveness of our DDCoT with existing pre-trained VLMs and multimodal reasoning models

We also validate the effectiveness of our DDCoT with existing pre-trained VLMs and multimodal reasoning models, as shown in Table 3. We observe that rationales generated by our proposed DDCoT is compatible with such pre-trained VLMs. Without correct rationales, existing pretrained

Model	NAT	SOC	LAN	TXT	IMG	NO	G1-6	G7-12	Avg
Flamingo [1]	40.28	46.13	29.25	47.45	40.08	33.66	39.35	40.67	39.83
Flamingo with our R	40.28	46.13	29.25	47.45	40.08	33.66	39.35	40.67	39.83
MiniGPT-4 [6]	40.28	46.13	29.25	47.45	40.08	33.66	39.35	40.67	39.83
MiniGPT-4 with our R	40.28	46.13	29.25	47.45	40.08	33.66	39.35	40.67	39.83

Table 3: The effectiveness of our DDCoT with Flamingo [1] and MiniGPT-4 [6]

VLMs and multimodal reasoning models have difficulty in complex reasoning tasks. Fortunately, the generalizable rationales generated by our DDCoT prompting can help existing VLMs to comprehend visual information and reason with rich knowledge, achieving significant improvement of 11.14% and 10.96% based on Flamingo [1] and Mini GPT-4 [6] as Table 3 shows.

	NoCaps			MSVD-QA
	CIDEr	SkipThoughtCS	EmbeddingAverageCS	GreedyMatchingScore
BLIP-2	76.15	49.84	89.20	77.94
Ours	46.26	84.78	92.35	79.12

Table 4: Addition experiments on NoCaps and MSVD-QA.

### A.5 Additional experiments on Captioning and Video Question Answering tasks

We extend our approach to more appropriate datasets. While other existing datasets may not fully exploit the benefits of our approach, we venture into exploring the captioning task on NoCaps and the video question-answering task on MSVD-QA, as shown in Table 4.

For captioning, we prompt the LLM to solve sub-problems derived from a simple caption, aiming to optimize and enrich it using the corresponding sub-answers. The substantial knowledge within LLM enables the generation of semantically enriched captions, leading to improvements in metrics evaluating sentence semantics, i.e. 34.94%, 3.15%, and 1.18% in terms of SkipThoughtCS, EmbeddingAverageCS, and GreedyMatchingScore. Note that the CIDEr metric to evaluate ours is limiting. It is designed to measure the similarity between the tested caption and reference captions without considering the diversity and high-level semantics.

For video question answering, LLM deconstructs problems like the decomposition step on ScienceQA. We sample video frames for VQA recognition and integrate frame information for multimodal rationale and answers. Leveraging the sequence understanding in LLM and visual information returned by the VQA model, we achieve a 4.9% improvement over BLIP-2.

Note that we randomly evaluated only 1000 images from NoCaps and 1000 videos from the MSVD test dataset in a zero-shot setting.

## B Human Evaluation

This section introduces the details of our human evaluation. Figure 4(b) shows an example of our question. Each sample comprises the question, context, options, and image. Evaluators are asked to rate the rationales generated by GPT [2], MM-COT [5], and our method in five aspects: relevance (related to the question), correctness (accuracy of reasoning and answer), completeness (logical reasoning’s comprehensiveness), coherence (consistency of reason), and explainability (interpretability of reasoning and answer). The rating scale ranges from 0 to 5. Additionally, the questions and rationales are organized into 12 groups, with each group assigned to three evaluators. Finally, we average the scores for each aspect of each rationale, resulting in overall scores and ratios relative to the maximum score.

## C Detailed Prompts

This section introduces more details about our zero-shot DDCoT prompting. We employ ChatGPT[4] as the most important component of our rationale generator. Specifically, it is utilized in breaking

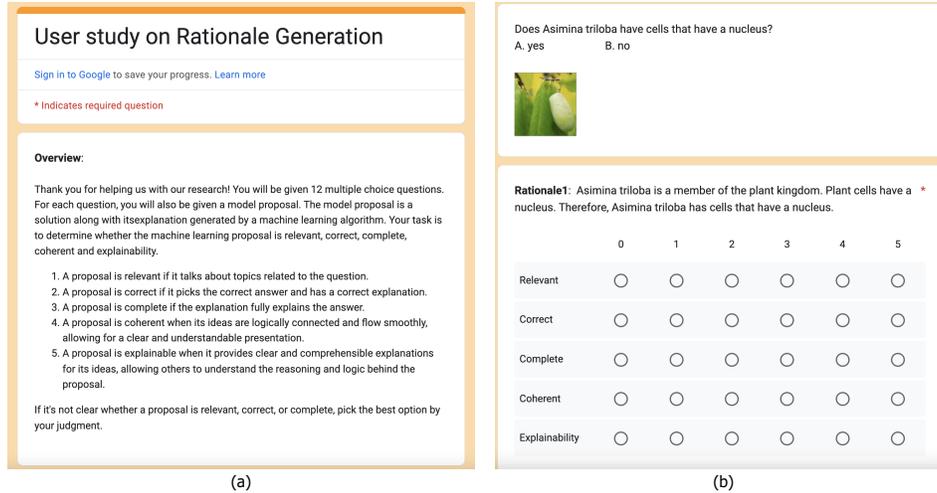


Figure 4: Interface of human evaluation. Figure(a) shows the instructions, figure(b) shows one example of our question.

duties of reasoning and recognition step and joint reasoning step. Figure 5 shows complete prompts for zero-shot DDCoT prompting.

## D Case Studies

To better understand the effectiveness of our proposed method in generating rationales, we randomly selected several cases from the test set along with the process of rationale generation. Figure 6 showcases several map-related questions, demonstrating how our method integrates simple visual features (such as the shape of highlighted areas) with common knowledge to obtain correct reasoning and answers. In the examples presented in Figure 7, our method successfully identifies within the images, acquiring relevant knowledge. Figure 8 illustrates four more complex questions, where our method leverages information obtained from the images to perform intricate reasoning. However, when it comes to the complex interaction between images and textual context, our method still fall into hallucinations, leading to erroneous reasoning and answers.

## E Limitations

While we have succeeded in mitigating a portion of the hallucination problem arising from multimodal inputs, it is important to acknowledge that this problem is not entirely resolved. As illustrated in Figure 8, our model is also susceptible to the risk of hallucinations. Investigating methods to suppress hallucinations is a potential topic for further research and exploration. In addition, we did not use extra image-text pairs to pre-train the alignment between vision and language modalities. Such pre-training is expected to further improve the alignment for joint reasoning.

As our approach involves zero-shot prompting of the LLM to generate rationales, there exists a potential risk of inheriting social biases from the LLM. These biases, which encompass cultural, ethical, and various other dimensions, might be reflected in the generated rationales, potentially leading to adverse effects on users. To mitigate this issue in the future, potential solutions could involve designing constraints at each prompting stage or utilizing more advanced LLMs trained on unbiased resources.

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.

DDCoT
<p><b>Breaking Reasoning Down to Recognition Steps with Negative-Space Prompting</b></p> <p><b>System:</b> You are a helpful, highly intelligent guided assistant. You will do your best to guide humans in choosing the right answer to the question. <u>Note that insufficient information to answer questions is common, because you do not have any information about the picture.</u> The final answer should be one of the options.</p> <p><b>User:</b> Given the context, questions and options, <u>please think step-by-step</u> about the preliminary knowledge to answer the question, <u>deconstruct the question</u> as completely as possible <u>down to</u> necessary <u>sub-questions</u> based on context, questions and options. Then with the aim of helping humans answer the original question, try to answer the sub-questions. The expected answering form is as follows:  Sub-questions:  1. &lt;sub-question 1&gt;  2. &lt;sub-question 2&gt;  ...  Sub-answers:  1. &lt;sub-answer 1&gt; or 'Uncertain'  2. &lt;sub-answer 2&gt; or 'Uncertain'  ...  Answer: &lt;One of the options&gt; or 'Uncertain'</p> <p>For a question, <u>assume that you do not have any information about the picture</u>, but <u>try to answer the sub-questions</u> and prioritize whether your general knowledge can answer it, and then consider whether the context can help. If sub-questions can be answered, then answer in as short a sentence as possible. <u>If sub-questions cannot be determined without information in images, please formulate corresponding sub-answer into "Uncertain".</u>  Only use "Uncertain" as an answer if it appears in the sub-answers. All answers are expected as concise as possible.  Here is an attempt:  Context: {context}  Has An Image: {has_image}  Question: {question}  Options: {option}"</p> <hr/> <p><b>Visual Recognition to Obtain Visual Complements</b></p> <p style="text-align: center;">~ Arbitrary Visual Question Answering (VQA) Model ~</p> <hr/> <p><b>Integrate to Joint Reasoning</b></p> <p><b>System:</b> You are a helpful, highly intelligent teacher. You will not only do your best to guide humans to the correct answer, but you will also give the rationales as a reference.</p> <p><b>User:</b> Given the context, questions, options, supplementary information, <u>think step by step</u> and <u>answer the questions</u>. Please note that we need not only the answer, but more importantly the rationales of getting the answer. The expected answering form is as follows:  Rationale: &lt;rationale&gt;  Answers: &lt;one of the options&gt;</p> <p>Please <u>note that the supplementary information given may not always be valid</u>. Please <u>select valid information to form the rationale</u> and choose the relatively correct option as your answer.  Here is an attempt:  Context: {context}  Has An Image: {has_image}  Question: {question}  Options: {option}  Supplementary information: {supplementary_information}</p>

Figure 5: The complete prompts of our zero-shot DDCoT prompting. The underlined parts are introduced in Section 3.2 of the submitted paper. The yellow background represents designs for critical thinking, and the blue background represents designs for differentiate duties of reasoning and recognition. “{ }” denotes the corresponding inputs, and “{supplementary\_information}” represents the sub-questions and sub-answers obtained before.

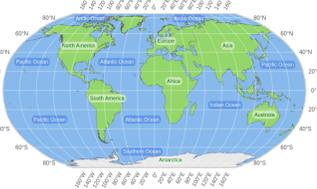
<p style="text-align: center;"><b>Question</b></p> <p>What is the capital of North Carolina?</p> <p>Context: Select the better answer.</p> <p>Options:</p> <p>A. Saint Paul B. Fargo C. Raleigh D. Charlotte</p> 	<p style="text-align: center;"><b>Ours Outputs</b></p> <p><b>Sub-question:</b>  (1) What is the location of North Carolina?  (2) What is the capital of North Carolina?</p> <p><b>Sub-answers:</b> (1) North Carolina is a state located in the southeastern region of the United States. (2) The capital of North Carolina is Raleigh.</p> <p><b>VQA:</b> Understanding[sub-questions, image]</p> <p><b>Sub-answers':</b> (1) North Carolina is a state located in the southeastern region of the United States. (2) <b>The capital of North Carolina is Raleigh.</b></p> <p><b>Rationale:</b> The question is asking for the capital of North Carolina. The supplementary information confirms that North Carolina is a state located in the southeastern region of the United States and the capital of North Carolina is Raleigh.</p>
<p style="text-align: center;"><b>Question</b></p> <p>Which property matches this object?</p> <p>Context: Select the better answer.</p> <p>Options:</p> <p>A. North America B. South America  C. Antarctica D. Australia</p> 	<p style="text-align: center;"><b>Ours Outputs</b></p> <p><b>Sub-question:</b>  (1) What is the shape of the highlighted area?  (2) What are the geographical features of the highlighted area?  (3) What are the surrounding bodies of water of the highlighted area?</p> <p><b>Sub-answers:</b> (1) Uncertain (2) Uncertain (3) Uncertain</p> <p><b>VQA:</b> Understanding[sub-questions, image]</p> <p><b>Sub-answers':</b> (1) a circle (2)a mountain range (3) atlantic ocean, pacific ocean, and arctic ocean</p> <p><b>Rationale:</b> <b>The highlighted area is a circle with a mountain range, and it is surrounded by the Atlantic Ocean, Pacific Ocean, and Arctic Ocean.</b> Based on this information, we can eliminate Antarctica and Australia as they are not surrounded by these three oceans. The mountain range in the highlighted area is likely the Rocky Mountains, which are located in <b>North America.</b></p>
<p style="text-align: center;"><b>Question</b></p> <p>Which country is highlighted?</p> <p>Context: Select the better answer.</p> <p>Options:</p> <p>A. Haiti B. Saint Lucia  C. Saint Vincent and the Grenadines</p> 	<p style="text-align: center;"><b>Ours Outputs</b></p> <p><b>Sub-question:</b>  (1) What is the image showing?  (2) What are the distinguishing features of the highlighted country?  (3) What are the geographical locations of the countries in the options?</p> <p><b>Sub-answers:</b> (1) Uncertain (2) Uncertain (3) Uncertain</p> <p><b>VQA:</b> Understanding[sub-questions, image]</p> <p><b>Sub-answers':</b> (1) the location of the caribbean (2) the southeastern coast is a peninsula (3) the caribbean islands are located in the western hemisphere</p> <p><b>Rationale:</b> <b>The image is showing the location of the Caribbean.</b> The supplementary information tells us that the highlighted country <b>has a southeastern coast that is a peninsula.</b> Based on the options, only Haiti and Saint Lucia have a southeastern coast that is a peninsula. However, Haiti does not have a peninsula that is as pronounced as the one in the image.</p>
<p style="text-align: center;"><b>Question</b></p> <p>Which of these continents does the prime meridian intersect?</p> <p>Context: Select the better answer.</p> <p>Options:</p> <p>A. South America B. Asia C. Antarctica</p> 	<p style="text-align: center;"><b>Ours Outputs</b></p> <p><b>Sub-question:</b>  (1) What is the prime meridian?  (2) Which continents does the prime meridian pass through?  (3) Is the prime meridian shown in the image?</p> <p><b>Sub-answers:</b> (1) The prime meridian is the line of longitude that is designated as 0° longitude. It is the line that divides the Earth into the Eastern and Western Hemispheres. (2)The prime meridian passes through Europe, Africa, and Antarctica. (3) Uncertain</p> <p><b>VQA:</b> Understanding[sub-questions, image]</p> <p><b>Sub-answers':</b> (1) The prime meridian is the line of longitude that is designated as 0° longitude. It is the line that divides the Earth into the Eastern and Western Hemispheres. (2)The prime meridian passes through <b>Europe, Africa, and Antarctica.</b> (3) no</p> <p><b>Rationale:</b> The prime meridian is the line of longitude that is designated as 0° longitude. It is the line that divides the Earth into the Eastern and Western Hemispheres. The supplementary information states that the prime meridian passes through Europe, Africa, and Antarctica. The image does not show the prime meridian.</p>
<p style="text-align: center;"><b>GPT3(COT) Outputs</b></p> <p><b>BECAUSE:</b> Raleigh is the capital of North Carolina.</p>	
<p style="text-align: center;"><b>GPT3(COT) Outputs</b></p> <p><b>BECAUSE:</b> Antarctica is the world's fifth largest continent. It is also the coldest, windiest, and driest continent.</p>	
<p style="text-align: center;"><b>GPT3(COT) Outputs</b></p> <p><b>BECAUSE:</b> The country highlighted is <b>Saint Vincent and the Grenadines.</b></p>	
<p style="text-align: center;"><b>GPT3(COT) Outputs</b></p> <p><b>BECAUSE:</b> The prime meridian intersects with <b>Asia.</b></p>	

Figure 6: Cases that need to understand the map.

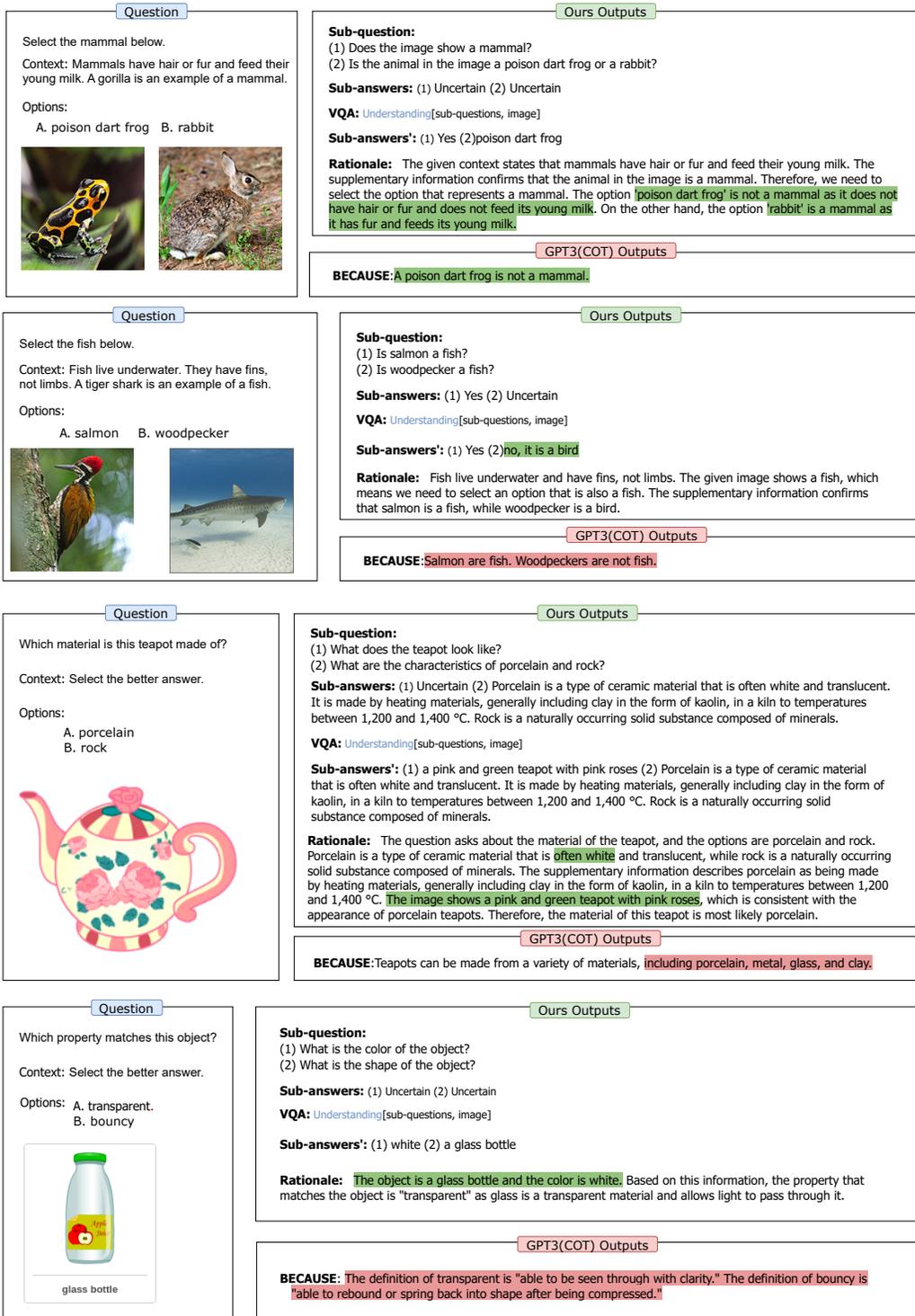


Figure 7: Cases that need to identify objects and acquire relevant knowledge.

**Question**

Compare the average kinetic energies of the particles in each sample. Which sample has the higher temperature?

Context: The diagrams below show two pure samples of gas in identical closed, rigid containers. Each colored ball represents one gas particle. Both samples have the same number of particles.

Options:  
A. sample B B. sample A  
C. neither; the samples have the same temperature

Sample A  
Mass of each particle: 44 u  
Average particle speed: 770 m/s

Sample B  
Mass of each particle: 44 u  
Average particle speed: 700 m/s

**Ours Outputs**

**Sub-question:**  
(1) What is the relationship between the average kinetic energy of gas particles and temperature?  
(2) Is there any difference in the speed or mass of the particles in each sample?  
(3) Is there any difference in the number of particles in each sample?

**Sub-answers:** (1) The average kinetic energy of gas particles is directly proportional to temperature. (2) It is uncertain whether there is any difference in the speed or mass of the particles in each sample. (3) Both samples have the same number of particles.

**VQA:** Understanding[sub-questions, image]

**Sub-answers':** (1) The average kinetic energy of gas particles is directly proportional to temperature. (2) the particles in sample a are moving at a much faster speed than the particles in sample b. (3) Both samples have the same number of particles.

**Rationale:** The average kinetic energy of gas particles is directly proportional to temperature. The particles in sample A are moving at a much faster speed than the particles in sample B. Both samples have the same number of particles. Therefore, sample A has a higher temperature than sample B.

**GPT3(COT) Outputs**

**BECAUSE:** The average kinetic energy of the particles in sample B is higher than the average kinetic energy of the particles in sample A. This means that sample B has a higher temperature.

**Question**

Which solution has a higher concentration of blue particles?

Context: The diagram below is a model of two solutions. Each blue ball represents one particle of solute.

Options:  
A. sample A B. sample B  
C. neither; their concentrations are the same

Solvent volume: 20 mL      Solvent volume: 20 mL

Solution A      Solution B

**Ours Outputs**

**Sub-question:**  
(1) What do the blue balls represent in the diagram?  
(2) Is there a difference in the number of blue balls in Solution A and Solution B?

**Sub-answers:** (1) The blue balls represent particles of solute. (2) Uncertain

**VQA:** Understanding[sub-questions, image]

**Sub-answers':** (1) The blue balls represent particles of solute. (2) no

**Rationale:** The blue balls represent particles of solute in both Solution A and Solution B. As there is no difference in the number of blue balls in both solutions, the concentration of blue particles in both solutions is the same.

**GPT3(COT) Outputs**

**BECAUSE:** The concentration of a solution is the amount of solute in a given amount of solvent. In other words, it is the amount of solute per unit volume of solvent. The concentration of a solution can be expressed in several ways. One way is to express it as the number of moles of solute per liter of solution. Another way is to express it as the number of grams of solute per liter of solution. The concentration of Solution A is 0.5 moles per liter. The concentration of Solution B is 1.0 moles per liter. Therefore, Solution B has a higher concentration of blue particles.

**Question**

Which month has the lowest average precipitation in London?

Context: Use the graph to answer the question below.

Options:  
A. July B. November C. February

Average precipitation in London, England

Precipitation (mm)

Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec

**Ours Outputs**

**Sub-question:**  
(1) What does the graph show?  
(2) How can we determine the average precipitation for each month?  
(3) Which month has the lowest average precipitation?

**Sub-answers:** (1) The graph shows the precipitation levels in London for each month of the year. (2) We can determine the average precipitation for each month by calculating the total precipitation for that month and dividing it by the number of years in the dataset. (3) The month with the lowest average precipitation in London is February.

**VQA:** Understanding[sub-questions, image]

**Sub-answers':** (1) The graph shows the precipitation levels in London for each month of the year. (2) We can determine the average precipitation for each month by calculating the total precipitation for that month and dividing it by the number of years in the dataset. (3) The month with the lowest average precipitation in London is February.

**Rationale:** We need to look at the graph to determine the month with the lowest average precipitation in London. From the graph, we can see that the bar for July is the shortest, indicating the lowest precipitation level. To confirm this, we can calculate the average precipitation for each month by dividing the total precipitation for each month by the number of years in the dataset. Doing so, we will find that February has the lowest average precipitation.

**GPT3(COT) Outputs**

**BECAUSE:** The graph shows the average monthly precipitation in London from January to December. The months are listed across the bottom of the graph. The precipitation is measured in millimeters and is shown by the blue bars. The tallest blue bar shows that the month with the highest average precipitation is July. The shortest blue bar shows that the month with the lowest average precipitation is February.

**Question**

What is the probability that a rat produced by this cross will have a dwarf body?

Context: In a group of rats, some individuals have a normal-sized body and others have a dwarf body. In this group, the gene for the body size trait has two alleles. The allele for a dwarf body (b) is recessive to the allele for a normal-sized body (B).

Options:  
A. 3/4  
B. 0/4  
C. 2/4  
D. 1/4  
E. 4/4

	B	b	
B	BB	Bb	
B	BB	Bb	

**Ours Outputs**

**Sub-question:**  
(1) What are the genotypes of the parent rats?  
(2) What are the possible genotypes of the offspring?  
(3) What is the probability of an offspring inheriting the recessive allele for dwarf body?

**Sub-answers:** (1) The Punnett square shows that the genotypes of the parent rats are Bb and Bb. (2) The possible genotypes of the offspring are BB, Bb, and bb. (3) The probability of an offspring inheriting the recessive allele for dwarf body (bb) is 1/4 or 25%.

**VQA:** Understanding[sub-questions, image]

**Sub-answers':** (1) The Punnett square shows that the genotypes of the parent rats are Bb and Bb. (2) The possible genotypes of the offspring are BB, Bb, and bb. (3) The probability of an offspring inheriting the recessive allele for dwarf body (bb) is 1/4 or 25%.

**Rationale:** The Punnett square shows that the genotypes of the parent rats are Bb and Bb. The possible genotypes of the offspring are BB, Bb, and bb. The allele for a dwarf body (b) is recessive to the allele for a normal-sized body (B). Therefore, only the bb genotype will result in a dwarf body. The probability of an offspring inheriting the recessive allele for dwarf body (bb) is 1/4 or 25%.

**GPT3(COT) Outputs**

**BECAUSE:** A Punnett square is a tool that shows all the possible combinations of alleles that can result from a genetic cross. In this Punnett square, the alleles for the body size trait are shown on the top and the left side of the square. The alleles for the body size trait are B and b. The allele for a normal-sized body (B) is dominant to the allele for a dwarf body (b). The allele for a normal-sized body (B) is represented by a capital letter. The allele for a dwarf body (b) is represented by a lowercase letter. The allele for a normal-sized body (B) is dominant to the allele for a dwarf body (b).

Figure 8: Complex and Challenging cases.

- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [3] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- [4] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [5] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.
- [6] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.