

# GPT4RoI: INSTRUCTION TUNING LARGE LANGUAGE MODEL ON REGION-OF-INTEREST

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Visual instruction tuning large language model (LLM) on image-text pairs has achieved general-purpose vision-language abilities. However, the lack of region-text pairs limits their advancements to fine-grained multimodal understanding. In this paper, we propose *spatial instruction tuning*, which introduces the reference to the region-of-interest (RoI) in the instruction. Before sending to LLM, the reference is replaced by RoI features and interleaved with language embeddings as a sequence. Our model GPT4RoI, trained on 7 region-text pair datasets, brings an unprecedented interactive and conversational experience compared to previous image-level models. (1) *Interaction beyond language*: Users can interact with our model by both language and drawing bounding boxes to flexibly adjust the referring granularity. (2) *Versatile multimodal abilities*: A variety of attribute information within each RoI can be mined by GPT4RoI, *e.g.*, color, shape, material, action, *etc.* Furthermore, it can reason about multiple RoIs based on common sense. On the Visual Commonsense Reasoning (VCR) dataset, GPT4RoI achieves a remarkable accuracy of 81.6%, surpassing all existing models by a significant margin (the second place is 75.6%) and almost reaching human-level performance of 85.0%. The code, dataset, and demo can be found at <https://github.com/Anonymous-Researcher1/GPT4RoI>.

## 1 INTRODUCTION

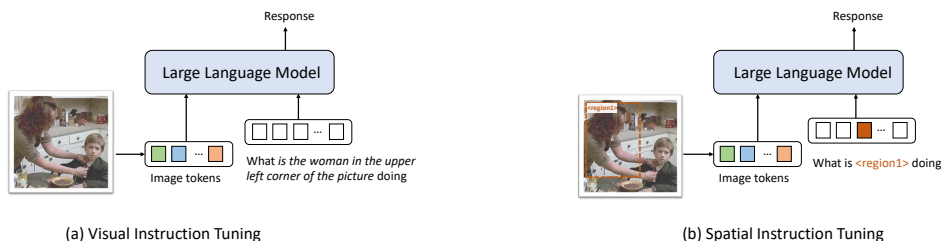


Figure 1: Comparison of visual instruction tuning on image-text pairs and spatial instruction tuning on region-text pairs. The bounding box and text description of each object are provided in region-text datasets. During training, the bounding box is from annotations, and in inference, it can be provided by user or any off-the-shelf object detector

Recent advancements of large language models (LLM) have shown incredible performance in solving natural language processing tasks in a human-like conversational manner, for example, commercial products (OpenAI, 2022; Anthropic, 2023; Google, 2023; OpenAI, 2023) and community open-source projects (Touvron et al., 2023a;b; Taori et al., 2023; Chiang et al., 2023; Du et al., 2022; Sun & Xipeng, 2022). Their unprecedented capabilities present a promising path toward general-purpose artificial intelligence models. Witnessing the power of LLM, the field of multimodal models (Yang et al., 2023c; Huang et al., 2023; Girdhar et al., 2023; Driess et al., 2023) is developing a new technology direction to leverage LLM as the universal interface to build general-purpose models,

Model	Image	Region	Multi-Region	Multi-Round Dialogue	End-to-End Model
Visual ChatGPT (Wu et al., 2023)	✓	✗	✗	✓	✗
MiniGPT-4 (Zhu et al., 2023)	✓	✗	✗	✓	✓
LLaVA (Liu et al., 2023a)	✓	✗	✗	✓	✓
InstructBLIP (Dai et al., 2023)	✓	✗	✗	✓	✓
MM-REACT (Yang et al., 2023c)	✓	✓	✓	✓	✗
InternGPT (Liu et al., 2023d)	✓	✓	✓	✓	✗
VisionLLM (Wang et al., 2023b)	✓	✓	✗	✗	✓
CaptionAnything (Wang et al., 2023a)	✓	✗	✗	✗	✗
DetGPT (Pi et al., 2023)	✓	✓	✗	✓	✗
GPT4RoI	✓	✓	✓	✓	✓

Table 1: Comparisons of vision-language models. Our GPT4RoI is an end-to-end model that supports region-level understanding and multi-round conversation.

where the feature space of a specific task is tuned to be aligned with the feature space of pre-trained language models.

As one of the representative tasks, vision-and-language models align the vision encoder feature to LLM by instruction tuning on image-text pairs, such as MiniGPT-4 (Zhu et al., 2023), LLaVA (Liu et al., 2023a), InstructBLIP (Dai et al., 2023), etc. Although these works achieve amazing multimodal abilities, their alignments are only on image-text pairs (Chen et al., 2015; Sharma et al., 2018; Changpinyo et al., 2021; Ordonez et al., 2011; Schuhmann et al., 2021), the lack of region-level alignment limits their advancements to more fine-grained understanding tasks such as region caption (Krishna et al., 2017) and reasoning (Zellers et al., 2019a). To enable region-level understanding in vision-language models, some works attempt to leverage external vision models, for example, MM-REACT (Yang et al., 2023c), InternGPT (Liu et al., 2023d) and DetGPT (Pi et al., 2023), as shown in Table 1. However, their non-end-to-end architecture is a sub-optimal choice for general-purpose multi-modal models.

Considering the limitations of previous works, our objective is to construct an end-to-end vision-language model that supports fine-grained understanding on region-of-interest. Since there is no operation that can refer to specific regions in current image-level vision-language models (Zhu et al., 2023; Liu et al., 2023a; Zhang et al., 2023c; Dai et al., 2023), our key design is to incorporate references to bounding boxes into language instructions, thereby upgrading them to the format of *spatial instructions*. For example, as shown in Figure 1, when the question is “*what is <region1> doing?*”, where the *<region1>* refers to a specific region-of-interest, the model will substitute the embedding of *<region1>* with the region feature extracted by the corresponding bounding box. The region feature extractor can be flexibly implemented by RoIAlign (He et al., 2017) or Deformable attention (Zhu et al., 2020).

To establish fine-grained alignment between vision and language, we involve region-text datasets in our training, where the bounding box and the text description of each region are provided. The datasets are consolidated from publicly available ones including COCO object detection (Lin et al., 2014), RefCOCO (Yu et al., 2016), RefCOCO+ (Yu et al., 2016), RefCOCOg (Mao et al., 2016), Flickr30K entities (Plummer et al., 2015), Visual Genome(VG) (Krishna et al., 2017) and Visual Commonsense Reasoning(VCR) (Zellers et al., 2019a). These datasets are transformed into spatial instruction tuning format. Moreover, we incorporate the LLaVA150K dataset (Liu et al., 2023a) into our training process by utilizing an off-the-shelf detector to generate bounding boxes. This enhances our model’s ability to engage in multi-round conversations and generate more human-like responses.

The collected datasets are categorized into two types based on the complexity of the text. First, the plain-text data contains object category and simple attribute information. It is used for pre-training the region feature extractor without impacting the LLM. Second, the complex-text data often contains complex concepts or requires common sense reasoning. We conduct end-to-end fine-tuning of the region feature extractor and LLM for these data.

Benefiting from spatial instruction tuning, our model brings a new interactive experience, where the user can express the question to the model with language and the reference to the region-of-interest. This leads to new capacities beyond image-level understanding, such as region caption and complex region reasoning. As a generalist, our model GPT4RoI also shows its strong region understanding ability on three popular benchmarks, including the region caption task on Visual Genome (Krishna et al., 2017), the region reasoning task on Visual-7W (Zhu et al., 2016) and Visual Commonsense Reasoning (Zellers et al., 2019a) (VCR). Especially noteworthy is the performance on the most challenging VCR dataset, where GPT4RoI achieves an impressive accuracy of 81.6%, 6 points ahead of the second-place and nearing the human-level performance benchmarked at 85.0%.

In summary, our work makes the following contributions:

- We introduce spatial instruction, combining language and the reference to region-of-interest into an interleave sequence, enabling accurate region referring and enhancing user interaction.
- By spatial instruction tuning LLM with massive region-text datasets, our model can follow user instructions to solve diverse region understanding tasks, such as region caption and reasoning.
- Our method, as a generalist, outperforms the previous state-of-the-art approach on a wide range of region understanding benchmarks.

## 2 RELATED WORK

### 2.1 LARGE LANGUAGE MODEL

The field of natural language processing (NLP) has achieved significant development by the high-capability large language model (LLM). The potential of LLM is first demonstrated by pioneering works such as BERT (Devlin et al., 2018) and GPT (Radford et al., 2018). Then scaling up progress is started and leads to a series of excellent works, for example, T5 (Raffel et al., 2020), GPT-3 (Brown et al., 2020), Flan-T5 (Chung et al., 2022), PaLM (Chowdhery et al., 2022), etc. With the growth of training data and model parameters, this scaling up progress brings to a phenomenal product, ChatGPT (OpenAI, 2022). By generative pre-trained LLM and instruction tuning (Ouyang et al., 2022) on human feedback, ChatGPT shows unprecedented performance on conversations with humans, reasoning and planning tasks (Mu et al., 2023; Yang et al., 2023a; Bubeck et al., 2023), etc.

### 2.2 LARGE VISION-LANGUAGE MODEL

To utilize high-performance LLM to build up vision-language models, LLM as task coordinator is proposed. Given the user instruction, LLM parses the instruction and calls various external vision models. Some representative works are Visual ChatGPT (Wu et al., 2023), ViperGPT (Surís et al., 2023), MM-REACT (Yang et al., 2023c), InternGPT (Liu et al., 2023d), VideoChat (Li et al., 2023b), etc. Although these models largely expand the scope of multimodal models, they depend on external vision models and these non-end-to-end architectures are not the optimal choice for multi-modal models. To obtain end-to-end vision-language models, instruction tuning LLM on image-text pairs is proposed to align visual features with LLM and accomplish multimodal tasks in a unified way, for example, Flamingo (Alayrac et al., 2022), MiniGPT-4 (Zhu et al., 2023), LLaVA (Liu et al., 2023a), LLaMa-Adapter (Zhang et al., 2023c), InstructBLIP (Dai et al., 2023), MM-GPT (Gong et al., 2023), VPGTrans (Zhang et al., 2023a), etc. These models achieve amazing image-level multimodal abilities, while several benchmarks such as LVLm-eHub (Xu et al., 2023) and MMBench (Liu et al., 2023c) find that these models still have performance bottlenecks when need to be under specific region reference. Our GPT4RoI follows the research line of visual instruction tuning and moves forward region-level multimodal understanding tasks such as region caption (Krishna et al., 2017) and reasoning (Zellers et al., 2019a).

### 2.3 REGION-LEVEL IMAGE UNDERSTANDING

For region-level understanding, it is a common practice in computer vision to identify potential regions of interest first and then do the understanding. Object detection (Ren et al., 2015; Carion et al., 2020; Zhu et al., 2020; Zang et al., 2023) tackles the search for potential regions, which are generally accompanied by a simple classification task to understand the region’s content. To expand

the object categories, (Kamath et al., 2021; Liu et al., 2023b; Zhou et al., 2022; Li\* et al., 2022) learn from natural language and achieve amazing open-vocabulary object recognition performance. Region captioning (Johnson et al., 2015; Yang et al., 2017; Wu et al., 2022) provides more descriptive language descriptions in a generative way. Scene graph generation (Li et al., 2017; Tang et al., 2018; Yang et al., 2022) analyzes the relationships between regions by the graph. The VCR (Zellers et al., 2019b) dataset presents many region-level reasoning cases and (Yu et al., 2021; Su et al., 2019; Li et al., 2019b; Yao et al., 2022) exhibit decent performance by correctly selecting the answers in the multiple-choice format. However, a general-purpose region understanding model has yet to emerge. In this paper, by harnessing the powerful large language model (Touvron et al., 2023a; Chiang et al., 2023), GPT4RoI uses a generative approach to handle all these tasks. Users can complete various region-level understanding tasks by freely asking questions.

### 2.4 USING TEXTUAL COORDINATES AS THE GROUNDING TOKEN.

We compare the design philosophy with methods using textual coordinates as the grounding token and provide a brief overview of concurrent works, all of which can be found in the appendix.

## 3 METHOD: GPT4RoI

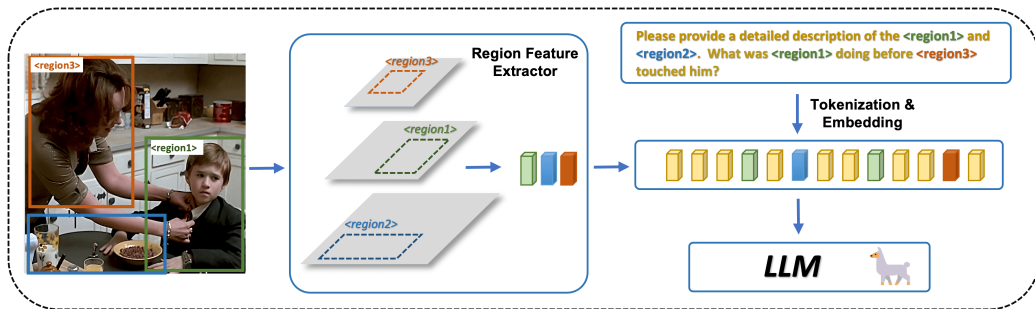


Figure 2: GPT4RoI is an end-to-end vision-language model for processing spatial instructions that contain references to the region-of-interest, such as  $\langle region\{i\} \rangle$ . During tokenization and conversion to embeddings, the embedding of  $\langle region\{i\} \rangle$  in the instruction is replaced with the RoIAlign results from multi-level image features. Subsequently, such an interleaved region feature and language embedding sequence can be sent to a large language model (LLM) for further processing. We also utilize the entire image feature to capture global information and omit it in the figure for brevity. A more detailed framework figure can be found in Figure 5 in the Appendix.

The overall framework of GPT4RoI consists of a vision encoder, a projector for image-level features, a region feature extractor, and a large language model (LLM). Compared to previous works (Zhu et al., 2023; Liu et al., 2023a), GPT4RoI stands out for its ability to convert instructions that include spatial positions into an interleaved sequence of region features and text embeddings, as shown in Figure 2.

### 3.1 MODEL ARCHITECTURE

We adopt the ViT-L/14 architecture from CLIP (Radford et al., 2021) as the vision encoder. Following (Liu et al., 2023a), we use the feature map of the penultimate transformer layer as the representation of the entire image, and then map the image feature embedding to the language space using a single linear layer as projector. Finally, we employ the Vicuna (Zheng et al., 2023), an instruction-tuned LLaMA (Touvron et al., 2023a), to perform further processing.

We utilize widely adopted modules in the field of object detection to construct our RoI feature extractor. To ensure a robust feature representation for regions of varying scales, we construct a multi-level image feature pyramid (Lin et al., 2017) by selecting four layers from the CLIP vision

encoder and fusing them with five lightweight scale shuffle modules (Zhang et al., 2023d). These layers are located at the second-to-last, fifth-to-last, eighth-to-last, and eleventh-to-last positions, respectively. Additionally, we incorporate feature coordinates (Liu et al., 2018a; Wang et al., 2020) for each level to address the problem of translation invariance in CNNs. This helps make the model sensitive to absolute position information, such as the description “*girl on left*” in Figure 3. Finally, we use RoIAlign to extract region-level features with an output size of  $14 \times 14$  (He et al., 2017), which ensures that sufficient detailed information is preserved. Moreover, all four level features are involved in the RoIAlign operation and fused into a single embedding as the representation of the region of interest (RoI) (Liu et al., 2018b).

### 3.2 TOKENIZATION AND EMBEDDING

To enable users to refer to regions of interest in text inputs, we define a special token  $\langle region\{i\} \rangle$ , which acts as the placeholder that will be replaced by the corresponding region feature after tokenization and embedding. One example is depicted in Figure 2. When a user presents a spatial instruction, “*What was  $\langle region1 \rangle$  doing before  $\langle region3 \rangle$  touched him?*”, the embedding of  $\langle region1 \rangle$  and  $\langle region3 \rangle$  are replaced by their corresponding region features. However, this replacement discards the references to different regions. To allow LLM to maintain the original references ( $region1$ ,  $region3$ ) in the response sequence, the instruction is modified to “*What was  $region1$   $\langle region1 \rangle$  doing before  $region3$   $\langle region3 \rangle$  touched him?*”. Then, LLM can generate a reply like “*The person in  $region1$  was eating breakfast before the person in  $region3$  touched them.*”

Regardless of the user instruction, we incorporate a prefix prompt, “*The  $\langle image \rangle$  provides an overview of the picture.*” The  $\langle image \rangle$  is a special token that acts as a placeholder, the embedding of which would be replaced by image features of the vision encoder. These features enable LLM to receive comprehensive image information and obtain a holistic understanding of the visual context.

### 3.3 SPATIAL INSTRUCTION TUNING

Our model is trained using a next-token prediction loss (Liu et al., 2023a; Zhu et al., 2023), where the model predicts the next token in a given input text sequence. The training details are in Section A.2 in the Appendix.

We transform annotations into instruction tuning format by creating a question that refers to the mentioned region for each region-text annotation. We partition the available region-text data into two groups, employing each in two distinct training stages. In the first stage, we attempt to align region features with word embeddings in language models using simple region-text pairs that contain color, position, or category information. The second stage is designed to handle more complex concepts, such as actions, relationships, and common sense reasoning. Furthermore, we provide diverse instructions for these datasets to simulate chat-like input in this stage.

**Stage 1: Pre-training** In this stage, we first load the weights of LLaVA (Liu et al., 2023a) after its initial stage of training, which includes a pre-trained vision encoder, a projector for image-level features, and an LLM. We only keep the region feature extractor trainable and aim to align region features with language embedding by collecting short text and bounding box pairs. These pairs are from both normal detection datasets and referring expression detection datasets, which have short expressions. The objective is to enable the model to recognize categories and simple attributes of the region in an image, which are typically represented by a short text annotation (usually within 5 words). Specifically, we utilize COCO (Lin et al., 2014), RefCOCO (Yu et al., 2016), and RefCOCO+ (Yu et al., 2016) datasets in this stage.

As shown in Table 2, for COCO detection data, we first explain the task in the prompt and then convert the annotations to a single-word region caption task. For RefCOCO and RefCOCO+, we also give task definitions first and train the model to generate descriptions containing basic attributes of the region. Only the description of the region (*in red color*) will be used to calculate the loss.

After this training stage, GPT4RoI can recognize categories, simple attributes, and positions of regions in images, as shown in Figure 3.

**Stage 2: End-to-end fine-tuning** In this stage, we only keep the vision encoder weights fixed and train the region feature extractor, image feature projector, and LLM weights. Our main focus is to


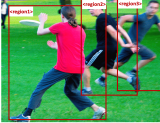
<p><b>Object Detection</b>                  In the conversation below, you simply answer the category name based on what you see in the imagery inside a particular region. I will give you only one region each time. Categories containing person, bicycle, car ...                  &lt;region1&gt; <b>person</b>                  &lt;region2&gt; <b>dog</b></p>	
<p><b>Referring Expression Comprehension</b>                  I will provide you with only one region containing only one object, although there may be other objects present in the image. It is recommended that you describe the object's relative position with respect to other objects in the image and its basic attributes.                  &lt;region1&gt; <b>red shirt girl</b>                  &lt;region2&gt; <b>guy in black</b>                  &lt;region3&gt; <b>right most person blurred</b></p>	

Table 2: The instruction template for Stage 1 training data: For both tasks, we begin by providing a description of the task definition and the expected answer. Then, we concatenate all region-text pairs into a sequence. For detection data, the format is <region{i}> category\_name. For referring expression comprehension, the format is <region{i}> description of region. Only the responses highlighted in red are used to calculate the loss.

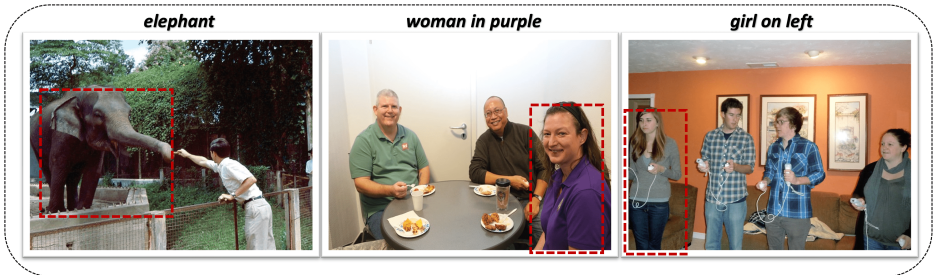


Figure 3: After Stage 1 training, GPT4RoI is capable of identifying the category of the region (elephant), simple attributes such as color (purple), and the position of the region (left).

enhance GPT4RoI’s ability to accurately follow user instructions and tackle complex single/multiple region understanding tasks. We tailor specific instructions for different tasks. For single region caption, we construct from Visual Genome (VG) region caption part (Krishna et al., 2017) and RefCOCOg (Mao et al., 2016). For multiple region caption, Flickr30k (Plummer et al., 2015) is converted to a multiple region caption task where the caption should include all visual elements emphasized by bounding boxes. To simulate user instruction, we create 20 questions for each caption task as shown in Table 8 and Table 9. For the region reasoning task, we modify Visual Commonsense Reasoning (VCR) (Zellers et al., 2019a) to meet the input format requirements and make it more similar to human input. The details of this process can be found in Section A.3.

To improve the capability of GPT4RoI for multi-round conversation and generate more human-like responses, we also involve the LLaVA150k (Liu et al., 2023a) visual instruction dataset in this stage. We employ an off-the-shelf LVIS detector (Fang et al., 2023) to extract up to 100 detection boxes per image. These boxes are then concatenated with the user instructions in the format “<region{i}> may feature a class\_name”. LLaVA150k significantly improves the capability of GPT4RoI for multi-round conversation .

After completing this training stage, GPT4RoI is capable of performing complex region understanding tasks based on user instructions, including region caption and reasoning, as demonstrated in Section 4.

#### 4 DEMONSTRATIONS

In this section, we compare the differences between the visual instruction tuning model LLaVA (Liu et al., 2023a) and our spatial instruction tuning model GPT4RoI. We demonstrate our new interactive approach and highlight its advanced capabilities in understanding multimodality.


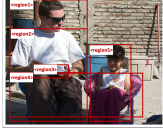

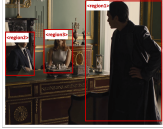
<p><b>Single Region Caption</b></p> <p>### Question: Can you provide me with a detailed description of the region in the picture marked by &lt;region1&gt;</p> <p>### Answer: <b>A man wearing a light blue t-shirt and jeans with his arms extend.</b></p>	
<p><b>Multiple Region Caption</b></p> <p>### Question: Can you provide a thorough description of the regions &lt;region1&gt;, &lt;region5&gt; in this image?</p> <p>### Answer: <b>A man is watching a girl draw as he eats</b></p> <p>### Question: Could you please give me a detailed description of areas &lt;region1&gt;, &lt;region2&gt;, &lt;region3&gt;, &lt;region4&gt;, &lt;region5&gt; ?</p> <p>### Answer: <b>A man in a white shirt with a plate of food sits outside in a folding chair with a little girl who is writing</b></p>	
<p><b>Single Region Reasoning</b></p> <p>### Question: How is &lt;region1&gt; feeling ?</p> <p>### Answer: <b>Person at region1 is feeling nervous</b></p> <p>### Question: Why do you think so?</p> <p>### Answer: <b>His hands are tied behind his back and he is surrounded by nazi guards. He is looking around to escape.</b></p>	
<p><b>Multi Region Reasoning</b></p> <p>### Question: Is &lt;region1&gt; happy to be speaking with &lt;region2&gt; and &lt;region3&gt; ?</p> <p>### Answer: <b>No, person at region1 is bothered by the conversation.</b></p> <p>### Question: What factors influenced your perspective?</p> <p>### Answer: <b>Person at region1 is standing with his hand on his hip in a defensive way.</b></p>	

Table 3: Instruction template for Stage 2 training data: During training, we randomly select one question for both single and multiple region caption tasks. For reasoning tasks, we modify the original questions to include a reference for each region so that GPT4RoI can mention them in its response. Only the response **in red color** and stop string **###** will be used to calculate the loss.

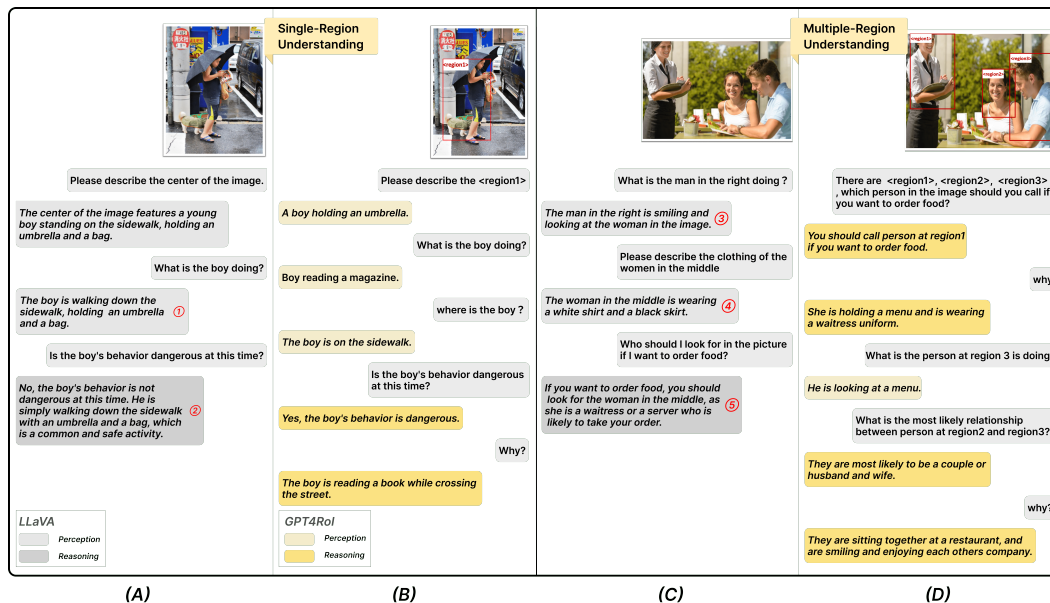


Figure 4: GPT4RoI and LLaVA dialogue performance showcase. Figures A and C demonstrate the dialogue scenarios of LLaVA when referring to a single instance and multiple instances solely using natural language in the conversation. On the other hand, Figures B and D showcase how GPT4RoI utilizes bounding boxes as references to address the same scenarios.

As shown in Figure 4.A, when we try to make LLaVA focus on the center region of the image, it only sees the boy holding an umbrella and a bag, but it misses the book. As a result, LLaVA gives a

wrong answer to the question “*What is the boy doing*” (Figure 4.A.①), and this leads to an incorrect conclusion that “*the boy’s behavior is not dangerous*” (Figure 4.A.②).

In comparison, as shown in Figure 4.B, our approach GPT4RoI efficiently recognizes visual details using the given bounding box. This allows it to accurately identify the action of “*reading a magazine.*” Furthermore, GPT4RoI demonstrates its reasoning abilities by correctly inferring that the “*boy’s behavior is dangerous*”, and giving a reasonable reason that “*the boy is reading a book while crossing the street.*”

When there are multiple instances in the image (as depicted in Figure 4.C), we attempt to refer to the corresponding instances as “*the right*” and “*the middle*”. However, LLaVA provides incorrect information by stating that the right man is “*looking at the women*” (as shown in Figure 4.C.③). Even more concerning, LLaVA overlooks the actual women in the middle and mistakenly associates the women on the left as the reference, resulting in completely inaccurate information (as shown in Figure 4.C.④ & ⑤).

In comparison, as shown in Figure 4.D, GPT4RoI is able to understand the user’s requirements, such as identifying the person to call when ordering food, and accurately recognize that the person in region1 fulfills this criterion. Additionally, it correctly recognizes that the person in region3 is “*looking at the menu*”. Importantly, GPT4RoI can also infer relationships between the provided regions based on visual observations. For example, it deduces that the likely relationship between region2 and region3 is that of a “*couple*”, providing a reasonable explanation that they “*are smiling and enjoying each other’s company.*”

## 5 QUANTITATIVE RESULTS

To quantitatively evaluate GPT4RoI, we have chosen three representative benchmarks to assess the region understanding capabilities. These benchmarks include the region caption task on Visual Genome (Krishna et al., 2017), the region reasoning task on Visual-7W (Zhu et al., 2016), and Visual Commonsense Reasoning (Zellers et al., 2019a) (VCR). In order to minimize the impact of specific dataset label styles and make evaluation metrics easier to calculate, we fine-tuned GPT4RoI on each benchmark using different task prompts. More details can be found in Section A.2 in the Appendix.

### 5.1 REGION CAPTION

We report the scores of BLEU, METEOR, ROUGE, and CIDEr for both GPT4RoI-7B and GPT4RoI-13B on the validation set of Visual Genome (Krishna et al., 2017). The grounding box in the annotation is combined with the task prompt in Appendix Table 7 to get the response.

Model	BLEU@4	METEOR	ROUGE	CIDEr
GRiT (Wu et al., 2022)	-	17.1	-	142.0
GPT4RoI-7B	11.5	17.4	35.0	145.2
GPT4RoI-13B	11.7	17.6	35.2	146.8

Table 4: Comparison of region caption ability on the validation dataset on Visual Genome. All methods employ ground truth bounding boxes and GPT4RoI can outperform previous state-of-the-art specialist GRiT.

The generalist approach GPT4RoI outperforms the previous state-of-the-art specialist model GRiT (Wu et al., 2022) by a significant margin, without any additional techniques or tricks. Additionally, we observe that the performance of GPT4RoI-7B and GPT4RoI-13B is comparable, suggesting that the bottleneck in performance lies in the design of the visual module and the availability of region-text pair data. These areas can be explored further in future work.

### 5.2 VISUAL-7W

Visual-7W (Zhu et al., 2016) is a PointQA dataset that contains a which box setting. Here, the model is required to choose the appropriate box among four options, based on a given description. For example, a question might ask, “*Which is the black machine under the sign?*”. This type of question



not only tests the model’s object recognition but also its ability to determine the relationship between objects.

To prevent information leakage, we remove overlapping images with the test set from Visual Genome (Krishna et al., 2017). The results clearly demonstrate that the 13B model outperforms the 7B model by a significant margin. This finding suggests that the reasoning ability heavily relies on the Large Language Model (LLM).

Model	LSTM-Att (Zhu et al., 2016)	CMNs (Hu et al., 2016)	12in1 (Lu et al., 2020)	GPT4RoI-7B	GPT4RoI-13B
Acc(%)	56.10	72.53	83.35	81.83	84.82

Table 5: Accuracy on Visual-7W test dataset.

### 5.3 VISUAL COMMONSENSE REASONING

Visual Commonsense Reasoning (VCR) offers a highly demanding scenario that necessitates advanced reasoning abilities, heavily relying on common sense. Given the question(Q), the model’s task is not only to select the correct answer(A) but also to select a rationale(R) that explains why the chosen answer is true. We give a more detailed explanation of each metric in our appendix

Model	Open Source	Parameters	Val Acc.(%)			Test Acc.(%)		
			Q → A	QA → R	Q → AR	Q → A	QA → R	Q → AR
ViLBERT (Lu et al., 2019)	Y	221M	72.4	74.5	54.0	73.3	74.6	54.8
Unicoder-VL (Li et al., 2019a)	Y	-	72.6	74.5	54.5	73.4	74.4	54.9
VLBERT-L (Su et al., 2019)	Y	383M	75.5	77.9	58.9	75.8	78.4	59.7
UNITER-L (Chen et al., 2020)	Y	303M	-	-	-	77.3	80.8	62.8
ERNIE-ViL-L (Yu et al., 2021)	Y	-	78.52	83.37	65.81	79.2	83.5	66.3
MERLOT (Zellers et al., 2021)	Y	223M	-	-	-	80.6	80.4	65.1
VILLA-L (Gan et al., 2020)	Y	-	78.45	82.57	65.18	78.9	82.8	65.7
RESERVE-L (Zellers et al., 2022)	Y	644M	-	-	-	84.0	84.9	72.0
VQA-GNN-L (Wang et al., 2022)	Y	1B+	-	-	-	85.2	86.6	74.0
GPT4RoI-7B	Y	7B+	<b>87.4</b>	<b>89.6</b>	<b>78.6</b>	-	-	-
VLUA+@Kuaishou	N	-	-	-	-	84.8	87.0	74.0
KS-MGSR@KDDI Research and SNAP	N	-	-	-	-	85.3	86.9	74.3
SP-VCR@Shopee	N	-	-	-	-	83.6	88.6	74.4
HunYuan-VCR@Tencent	N	-	-	-	-	85.8	88.0	75.6
Human Performance (Zellers et al., 2019a)	-	-	-	-	-	91.0	93.0	85.0
GPT4RoI-13B	Y	13B+	-	-	-	<b>89.4</b>	<b>91.0</b>	<b>81.6</b>

Table 6: Accuracy scores on VCR. GPT4RoI achieves state-of-the-art accuracy among all methods.

GPT4RoI shows significant improvements over the previous methods across all  $Q \rightarrow A$ ,  $QA \rightarrow R$ , and  $Q \rightarrow AR$  tasks. Notably, in the crucial  $Q \rightarrow AR$  task, GPT4RoI-13B achieves a performance of 81.6 accuracy, surpassing preceding methods by over 6 points, even outperforming confidential company-level results, which may take advantage of private data. Our totally open-source pipeline can make GPT4RoI a solid baseline. More importantly, this performance is almost reaching human-level performance of 85.0 accuracy, which shows that the multimodal ability of GPT4RoI is promising to be further developed to human intelligence. Furthermore, comparing GPT4RoI to previous methods, particularly observing the size of the language model used, also demonstrates the significant benefits of the Large Language Model (LLM) for visual reasoning tasks.

## 6 CONCLUSIONS

In this paper, we present GPT4RoI, an end-to-end vision-language model that can execute user instructions to achieve region-level image understanding. Our approach employs spatial instruction tuning for the large language model (LLM), where we convert the reference to bounding boxes from user instructions into region features. These region features, along with language embeddings, are combined to create an input sequence for the large language model. By utilizing existing open-source region-text pair datasets, we show that GPT4RoI enhances user interaction by accurately referring to regions and achieves impressive performance in region-level image understanding tasks.

## REFERENCES

- 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. 3
- Anonymous. Ins-detCLIP: Aligning detection model to follow human-language instruction. In *Submitted to The Twelfth International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=M0MF4t3hE9>. under review. 18
- Anthropic. Claude. <https://www.anthropic.com/index/introducing-claude>, 2023. 1
- 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
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023. 3
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pp. 213–229. Springer, 2020. 3
- Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3558–3568, 2021. 2
- Chi Chen, Ruoyu Qin, Fuwen Luo, Xiaoyue Mi, Peng Li, Maosong Sun, and Yang Liu. Position-enhanced visual instruction tuning for multimodal large language models. *arXiv preprint arXiv:2308.13437*, 2023a. 18
- Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023b. 18
- Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. *arXiv preprint arXiv:2109.10852*, 2021. 18
- Ting Chen, Saurabh Saxena, Lala Li, Tsung-Yi Lin, David J Fleet, and Geoffrey E Hinton. A unified sequence interface for vision tasks. *Advances in Neural Information Processing Systems*, 35: 31333–31346, 2022. 18
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 2
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning, 2020. 9
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>. 1, 4
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 3
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022. 3

- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. 2, 3
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 1
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 320–335, 2022. 1
- Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva-02: A visual representation for neon genesis. *arXiv preprint arXiv:2303.11331*, 2023. 6
- Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning, 2020. 9
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15180–15190, 2023. 1
- Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. Multimodal-gpt: A vision and language model for dialogue with humans, 2023. 3
- Google. Bard. <https://bard.google.com/>, 2023. 1
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017. 2, 5
- Ronghang Hu, Marcus Rohrbach, Jacob Andreas, Trevor Darrell, and Kate Saenko. Modeling relationships in referential expressions with compositional modular networks, 2016. 9
- Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning perception with language models. *arXiv preprint arXiv:2302.14045*, 2023. 1
- Justin Johnson, Andrej Karpathy, and Li Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning, 2015. 4
- Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr – modulated detection for end-to-end multi-modal understanding, 2021. 4
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017. 2, 3, 6, 8, 9, 16
- Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang, Lionel M Ni, and Heung-Yeung Shum. Mask dino: Towards a unified transformer-based framework for object detection and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3041–3050, 2023a. 18
- Gen Li, Nan Duan, Yuejian Fang, Ming Gong, Daxin Jiang, and Ming Zhou. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training, 2019a. 9
- KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhui Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023b. 3

- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language, 2019b. 4
- Liunian Harold Li\*, Pengchuan Zhang\*, Haotian Zhang\*, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded language-image pre-training. In *CVPR*, 2022. 4
- Yikang Li, Wanli Ouyang, Bolei Zhou, Kun Wang, and Xiaogang Wang. Scene graph generation from objects, phrases and region captions, 2017. 4
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014. 2, 5
- Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2117–2125, 2017. 4
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023a. 2, 3, 4, 5, 6
- Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution, 2018a. 5
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying dino with grounded pre-training for open-set object detection, 2023b. 4
- Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8759–8768, 2018b. 5
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023c. 3
- Zhaoyang Liu, Yinan He, Wenhai Wang, Weiyun Wang, Yi Wang, Shoufa Chen, Qinglong Zhang, Yang Yang, Qingyun Li, Jiashuo Yu, et al. Internchat: Solving vision-centric tasks by interacting with chatbots beyond language. *arXiv preprint arXiv:2305.05662*, 2023d. 2, 3
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks, 2019. 9
- Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 12-in-1: Multi-task vision and language representation learning, 2020. 9
- Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 11–20, 2016. 2, 6
- Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of thought. *arXiv preprint arXiv:2305.15021*, 2023. 3
- OpenAI. Chatgpt. <https://openai.com/blog/chatgpt>, 2022. 1, 3
- OpenAI. Gpt-4 technical report, 2023. 1
- Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24, 2011. 2

- 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. 3
- Renjie Pi, Jiahui Gao, Shizhe Diao, Rui Pan, Hanze Dong, Jipeng Zhang, Lewei Yao, Jianhua Han, Hang Xu, and Lingpeng Kong Tong Zhang. Detgpt: Detect what you need via reasoning. *arXiv preprint arXiv:2305.14167*, 2023. 2
- Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pp. 2641–2649, 2015. 2, 6
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. *OpenAI*, 2018. 3
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021. 4
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 3
- Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Erix Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. *arXiv preprint arXiv:2311.03356*, 2023. 18
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28: 91–99, 2015. 3, 18
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 2
- Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8429–8438, 2019. doi: 10.1109/ICCV.2019.00852. 18
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernamed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2556–2565, 2018. 2
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019. 4, 9
- Tianxiang Sun and Qiu Xipeng. Moss. <https://github.com/OpenLMLab/MOSS>, 2022. 1
- Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. *arXiv preprint arXiv:2303.08128*, 2023. 3
- Kaihua Tang, Hanwang Zhang, Baoyuan Wu, Wenhan Luo, and Wei Liu. Learning to compose dynamic tree structures for visual contexts, 2018. 4
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023. 1

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a. 1, 4
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b. 1
- Teng Wang, Jinrui Zhang, Junjie Fei, Yixiao Ge, Hao Zheng, Yunlong Tang, Zhe Li, Mingqi Gao, Shanshan Zhao, Ying Shan, et al. Caption anything: Interactive image description with diverse multimodal controls. *arXiv preprint arXiv:2305.02677*, 2023a. 2
- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023b. 2
- Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pp. 649–665. Springer, 2020. 5
- Yanan Wang, Michihiro Yasunaga, Hongyu Ren, Shinya Wada, and Jure Leskovec. Vqa-gnn: Reasoning with multimodal semantic graph for visual question answering, 2022. 9
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023. 2, 3
- Jialian Wu, Jianfeng Wang, Zhengyuan Yang, Zhe Gan, Zicheng Liu, Junsong Yuan, and Lijuan Wang. Grit: A generative region-to-text transformer for object understanding, 2022. 4, 8
- Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models, 2023. 3
- Jiange Yang, Wenhui Tan, Chuhao Jin, Bei Liu, Jianlong Fu, Ruihua Song, and Limin Wang. Pave the way to grasp anything: Transferring foundation models for universal pick-place robots. *arXiv preprint arXiv:2306.05716*, 2023a. 3
- Jingkang Yang, Yi Zhe Ang, Zujin Guo, Kaiyang Zhou, Wayne Zhang, and Ziwei Liu. Panoptic scene graph generation, 2022. 4
- Linjie Yang, Kevin Tang, Jianchao Yang, and Li-Jia Li. Dense captioning with joint inference and visual context. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul 2017. 4
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9, 2023b. 18
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023c. 1, 2, 3
- Yuan Yao, Qianyu Chen, Ao Zhang, Wei Ji, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. Pevl: Position-enhanced pre-training and prompt tuning for vision-language models. *arXiv preprint arXiv:2205.11169*, 2022. 4
- Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*, 2023. 18
- Fei Yu, Jiji Tang, Weichong Yin, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. Ernie-vil: Knowledge enhanced vision-language representations through scene graph, 2021. 4, 9

- Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pp. 69–85. Springer, 2016. [2](#), [5](#)
- Yuhang Zang, Wei Li, Jun Han, Kaiyang Zhou, and Chen Change Loy. Contextual object detection with multimodal large language models, 2023. [3](#)
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6720–6731, 2019a. [2](#), [3](#), [6](#), [8](#), [9](#), [16](#)
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019b. [4](#), [17](#)
- Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models. In *Advances in Neural Information Processing Systems 34*, 2021. [9](#)
- Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yanpeng Zhao, Mohammadreza Salehi, Aditya Kusupati, Jack Hessel, Ali Farhadi, and Yejin Choi. Merlot reserve: Multimodal neural script knowledge through vision and language and sound. In *CVPR*, 2022. [9](#)
- Ao Zhang, Hao Fei, Yuan Yao, Wei Ji, Li Li, Zhiyuan Liu, and Tat-Seng Chua. Transfer visual prompt generator across llms. *arXiv preprint arXiv:2305.01278*, 2023a. [3](#)
- Ao Zhang, Liming Zhao, Chen-Wei Xie, Yun Zheng, Wei Ji, and Tat-Seng Chua. Next-chat: An lmm for chat, detection and segmentation. *arXiv preprint arXiv:2311.04498*, 2023b. [18](#)
- Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. [18](#)
- Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023c. [2](#), [3](#)
- Shilong Zhang, Xinjiang Wang, Jiaqi Wang, Jiangmiao Pang, Chengqi Lyu, Wenwei Zhang, Ping Luo, and Kai Chen. Dense distinct query for end-to-end object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7329–7338, June 2023d. [5](#), [18](#)
- Liang Zhao, En Yu, Zheng Ge, Jinrong Yang, Haoran Wei, Hongyu Zhou, Jianjian Sun, Yuang Peng, Runpei Dong, Chunrui Han, et al. Chatspot: Bootstrapping multimodal llms via precise referring instruction tuning. *arXiv preprint arXiv:2307.09474*, 2023. [18](#)
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. [4](#)
- Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022. [4](#)
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. [2](#), [3](#), [4](#), [5](#)
- Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. [2](#), [3](#)
- Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images, 2016. [3](#), [8](#), [9](#), [16](#)

## A APPENDIX

In this appendix, we provide a detailed method architecture figure. We then discuss training-related details, including hyperparameters and instruction templates used in each stage and task. Specifically, we [give an introduction for VCR dataset](#) and describe how we utilize the VCR dataset. We also [compare the design philosophy with methods using textual coordinates in LLM](#) and provide a [brief overview of concurrent works](#). Finally, we analyze some error cases and propose potential improvements for future exploration.

### A.1 DETAILED ARCHITECTURE

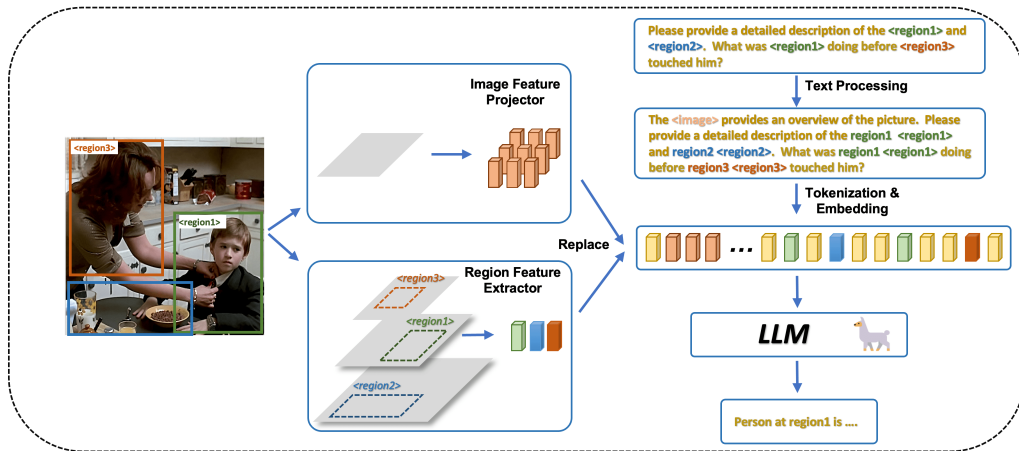


Figure 5: A more detailed framework of GPT4RoI.

Here is a more detailed framework of our approach, GPT4RoI.

1. We preprocess the input text by adding prefixes to retain both image information and pure text references for each region.
2. Next, we tokenize and embed the text. The image feature and region touched will replace the placeholders `<image>` and `<region[i]>` respectively.
3. The resulting interleaved sequence of region & image features and language embeddings is then fed into a large language model (LLM) for further processing.

### A.2 TRAINING DETAILS

**Dialogue model** The dialogue model in the demo is trained on 8 GPUs, each with 80G of memory. During the first training stage, a learning rate of  $2e-5$  is used with a cosine learning schedule. The batch size is 16 for 2 epochs, with a warm-up iteration set to 3000 and a warm-up ratio of 0.003. The weight decay for all modules was set to 0. During the second training stage, the learning rate is reduced to  $2e-5$  and the model is trained for 1 epoch. To enable end-to-end fine-tuning of the model, which includes a 7B Vicuna, Fully Sharded Data Parallel (FSDP) is enabled in PyTorch to save memory.

**Downstream tasks** We finetune on three datasets with different learning schedules and task prompts (as shown in Table 7). For the region caption task on Visual Genome (Krishna et al., 2017), we perform fine-tuning for 4 epochs with a learning rate of  $2e-5$ . As for Visual-7W (Zhu et al., 2016), we observe that it requires a smaller learning rate of  $1e-6$  to stabilize the training, which is also trained in 2 epochs. On the Visual Commonsense Reasoning (Zellers et al., 2019a), we fine-tune the model for 1 epoch using a learning rate of  $2e-5$ .

**Instruction of three downstream tasks.** The instructions for three downstream tasks are provided in Table 7.



**Region Caption Task on Visual Genome**

### Question: Can you give a description of the region mentioned by <region>

### Answer: **A man wearing a light blue t-shirt and jeans with his arms extended**

**Region Reasoning Task on Visual-7W**

### Question: <region1>,<region2>,<region3>,<region4> refers to specific areas within the photo along with their respective identifiers. I need you to answer the question. Questions are multiple-choice; you only need to pick the correct answer from the given options (A), (B), (C), or (D). Which is the black machine under the sign?

### Answer: **(A)**

**Region Reasoning Task on VCR**

**Q → A**

### Question: <region1>,<region2>,<region3>... refers to specific areas within the photo along with their respective identifiers. I need you to answer the question. Questions are multiple-choice; you only need to pick the correct answer from the given options (A), (B), (C), or (D).

How is 1 feeling ?

- (A),1 is feeling amused .
- (B),1 is upset and disgusted .
- (C),1 is feeling very scared .
- (D),1 is feeling uncomfortable with 3

### Answer: **(C)**

**QA → R**

### Question: <region1>,<region2>,<region3>... refers to specific areas within the photo along with their respective identifiers. I give you a question and its answer, I need you to provide a rationale explaining why the answer is right. Both questions are multiple-choice; you only need to pick the correct answer from the given options (A), (B), (C), or (D).

"How is 1 feeling ?" The answer is "1 is feeling very scared." What's the rationale for this decision?

- (A),1's face has wide eyes and an open mouth .
- (B),When people have their mouth back like that and their eyebrows lowered they are usually disgusted by what they see .
- (C),3,2,1 are seated at a dining table where food would be served to them . people unaccustomed to odd or foreign dishes may make disgusted looks at the thought of eating it .
- (D),1's expression is twisted in disgust .

### Answer: **(A)**

Table 7: Task prompt of three downstream tasks.

**Instruction of Single-Region Caption** The instructions for single-region caption are provided in Table 8. We randomly select one as the question in training.

**Instruction of Multi-Region Caption** The instructions for multi-region caption are provided in Table 9. We randomly select one as the question in training.

### A.3 VCR

**Introduction to the VCR Dataset** The Visual Commonsense Reasoning(VCR) dataset (Zellers et al., 2019b), comprises 290,000 multiple-choice questions obtained from 110,000 movie scenes. Each image in the dataset is annotated with a question that requires common-sense reasoning, along with its corresponding answer and the explanation for the answer. VCR is a particularly challenging dataset for comprehension and reasoning. It has gained attention from several well-known organizations, who have submitted their solutions on the leaderboard. The dataset's distinctive challenge is that a model not only needs to answer complex visual questions but also provide a rationale for why its answer is correct. The VCR task consists of two sub-tasks: Question Answering (Q→A) and Answer Justification (QA→R). In the Q→A setup, a model is given a question and must

select the correct answer from four choices. In the QA->R setup, a model is provided with a question and the correct answer, and it needs to justify the answer by selecting the most appropriate rationale from four choices. The performance of models is evaluated using the Q->AR metric, where accuracy is measured as the percentage of correctly answered questions along with the correct rationale.

**Preprocess of VCR** To construct a sequence of questions, we convert the explanation to a follow-up question and format them into a two-round conversation. Table 10 shows an example of the follow-up question that asks for the reasoning behind the answer.

The VCR dataset is valued for its diverse question-answer pairs that require referencing from prior question-answers to perform reasoning. Therefore, it’s crucial to assign a reference to each region in the dataset. We accomplish this by starting each conversation with a reference to all regions, e.g., *There are <region1>, <region2>... in the image.* This approach explicitly references every region, avoiding confusion in future analyses. Additionally, we substitute the corresponding *<region{i}>* in the answer with *category\_name at region{i}* to ensure a plain text output sequence.

#### A.4 TEXTUAL COORDINATES AS THE GROUNDING TOKEN

The key distinction lies in whether to incorporate the detection function into the LLM. For the method that uses textual coordinates as the grounding token, they have to solve the following challenge:

Aligning a large number of position tokens with their corresponding positions in the image by training on a large set of datasets. But this is actually a simple rule that can be naturally implemented with the operation in detection architectures.

Modeling geometric properties can be challenging. For example, if the ground truth box is  $\langle x_1 = 0, y_1 = 0, x_2 = 5, y_2 = 5 \rangle$ , a predicted box of  $\langle x_1 = 1, y_1 = 1, x_2 = 4, y_2 = 4 \rangle$  would be considered a better result than  $\langle x_1 = 1, y_1 = 1, x_2 = 8, y_2 = 8 \rangle$ . because it has a higher overlap with the ground truth. However, incorporating this geometric property into the next token prediction task using cross-entropy loss can be challenging. On the other hand, utilizing traditional loss functions such as L1 or IoU loss can naturally handle this geometric constraint.

Dense to Sparse (Ren et al., 2015; Zhang et al., 2023d) is a crucial design for detection performance, but embedding such an idea into the sequential form of LLM is challenging. We provide two pieces of evidence to support our argument

1. The performance of pix2seq (Chen et al., 2021; 2022), which utilizes object365 (Shao et al., 2019) pretrain, falls significantly behind the corresponding specialist (Zhang et al., 2022; Li et al., 2023a).
2. Even with scaled-up data and parameters, GPT4V still faces challenges in object counting (Yang et al., 2023b). However, this is a trivial task for detection methods.

Another approach is to use an external detector to find the potential region of interest, whereas LLM only focuses on analyzing the corresponding region of interest. This is the motivation of GPT4RoI. It requires much less data and allows for quick adaptation to specific domain problems with the corresponding detector. However, the drawback is that the framework may appear less elegant and it assumes input contains all regions of interest that need to be analyzed.

Both approaches have their advantages and disadvantages, and academic research in both directions is thriving (including concurrent works or follow-ups on GPT4RoI). For the first approach, relevant references include (Zhao et al., 2023; Chen et al., 2023b), while for the second approach, there are (Anonymous, 2023; Chen et al., 2023a) besides GPT4RoI. Additionally, there has been research that explores a fusion of the two approaches, as shown in references (You et al., 2023; Rasheed et al., 2023; Zhang et al., 2023b).

#### A.5 FAILURE CASE ANALYSIS

Due to limited data and instructions, GPT4RoI may fail in several landmark scenarios. We have conducted a thorough analysis and look forward to improving these limitations in future versions.

**Instruction obfuscation** As shown in Figure 6.(a), our multiple-region reasoning capability mainly relies on VCR, where we often use sentences that declare *<region1>*, *<region2>*, etc. at the beginning of the question. However, when users adopt the less common sentence structure to refer to regions, it

can often be confused with region captions that have the highest proportion in the dataset. As shown in Figure 6.(b), because our data and instructions are mainly generated by rules, our training data does not include content with the "respectively" instruction in multi-region scenarios. This can be resolved by adding specific instructions. In future versions, we aim to develop more diverse instructions, while ensuring data balance.



Figure 6: **GPT4RoI on instruction obfuscation.**

**Misidentification of fine-grained information within in region** Although GPT4RoI has improved the fine-grained perception ability of images compared to image-level vision language models, the limited amount of region-level data results in insufficient fine-grained alignment within regions. For example, in Figure 7.(a), the model incorrectly identifies the color of the helmet, and in Figure 7.(b), it misidentifies the object in the girl’s hand. Both cases generate the corresponding answers based on the most prominent feature within the region. Using semi-supervised methods to create more region-level data may address this issue.

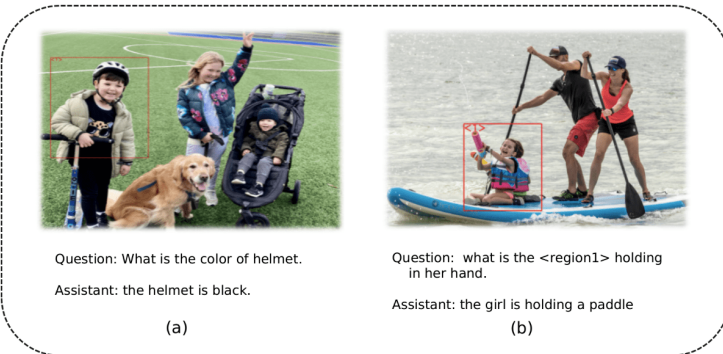


Figure 7: **GPT4RoI on Misidentification of fine-grained information.**

A.6 DISCUSSION

In our exploration, we find GPT4RoI produces failure cases as shown in Section. A.5. To further improve the performance, we identify the following potential directions:

- Model architecture. We find that 224 × 224 input image resolution struggles with understanding smaller regions. However, if we switch to a larger resolution, we must consider the potential burden on inference speed from global attention ViT architecture, while the more efficient CNN architecture or sliding window attention has no available pre-trained large-scale vision encoder like CLIP ViT-H/14.
- More region-text pair data. The amount of available region-text pairs is notably smaller than that of image-text pairs, which makes it challenging to sufficiently align region-level features

with language models. To tackle this issue, we may try to generate region-level pseudo labels by leveraging off-the-shelf detectors to generate bounding boxes for image-text data.

- **Region-level instructions.** Although we have generated instructions for each task from existing open-source datasets, users in practical applications may ask various questions about an arbitrary number of regions, and the existing data may not contain satisfactory answers. To tackle this issue, we suggest generating a new batch of spatial instructions through manual labeling or by leveraging ChatGPT or GPT4.
- **Interaction mode.** Currently, GPT4RoI only supports natural language and bounding box interaction. Incorporating more open-ended interaction modes such as point, scribble, or image-based search could further improve the user interaction experience.

1. Can you provide me with a detailed description of the region in the picture marked by <region1>?
2. I'm curious about the region represented by <region1> in the picture. Could you describe it in detail?
3. What can you tell me about the region indicated by <region1> in the image?
4. I'd like to know more about the area in the photo labeled <region1>. Can you give me a detailed description?
5. Could you describe the region shown as <region1> in the picture in great detail?
6. What details can you give me about the region outlined by <region1> in the photo?
7. Please provide me with a comprehensive description of the region marked with <region1> in the image.
8. Can you give me a detailed account of the region labeled as <region1> in the picture?
9. I'm interested in learning more about the region represented by <region1> in the photo. Can you describe it in detail?
10. What is the region outlined by <region1> in the picture like? Could you give me a detailed description, please?
11. Can you provide me with a detailed description of the region in the picture marked by <region1>, please?
12. I'm curious about the region represented by <region1> in the picture. Could you describe it in detail, please?
13. What can you tell me about the region indicated by <region1> in the image, exactly?
14. I'd like to know more about the area in the photo labeled <region1>, please. Can you give me a detailed description?
15. Could you describe the region shown as <region1> in the picture in great detail, please?
16. What details can you give me about the region outlined by <region1> in the photo, please?
17. Please provide me with a comprehensive description of the region marked with <region1> in the image, please.
18. Can you give me a detailed account of the region labeled as <region1> in the picture, please?
19. I'm interested in learning more about the region represented by <region1> in the photo. Can you describe it in detail, please?
20. What is the region outlined by <region1> in the picture like, please? Could you give me a detailed description?

Table 8: A list of instructions for single-region caption.

1. Could you please give me a detailed description of these areas [<region1>, <region2>, ...]?
2. Can you provide a thorough description of the regions [<region1>, <region2>, ...] in this image?
3. Please describe in detail the contents of the boxed areas [<region1>, <region2>, ...].
4. Could you give a comprehensive explanation of what can be found within [<region1>, <region2>, ...] in the picture?
5. Could you give me an elaborate explanation of the [<region1>, <region2>, ...] regions in this picture?
6. Can you provide a comprehensive description of the areas identified by [<region1>, <region2>, ...] in this photo?
7. Help me understand the specific locations labeled [<region1>, <region2>, ...] in this picture in detail, please.
8. What is the detailed information about the areas marked by [<region1>, <region2>, ...] in this image?
9. Could you provide me with a detailed analysis of the regions designated [<region1>, <region2>, ...] in this photo?
10. What are the specific features of the areas marked [<region1>, <region2>, ...] in this picture that you can describe in detail?
11. Could you elaborate on the regions identified by [<region1>, <region2>, ...] in this image?
12. What can you tell me about the areas labeled [<region1>, <region2>, ...] in this picture?
13. Can you provide a thorough analysis of the specific locations designated [<region1>, <region2>, ...] in this photo?
14. I am interested in learning more about the regions marked [<region1>, <region2>, ...] in this image. Can you provide me with more information?
15. Could you please provide a detailed description of the areas identified by [<region1>, <region2>, ...] in this photo?
16. What is the significance of the regions labeled [<region1>, <region2>, ...] in this picture?
17. I would like to know more about the specific locations designated [<region1>, <region2>, ...] in this image. Can you provide me with more information?
18. Can you provide a detailed breakdown of the regions marked [<region1>, <region2>, ...] in this photo?
19. What specific features can you tell me about the areas identified by [<region1>, <region2>, ...] in this picture?
20. Could you please provide a comprehensive explanation of the locations labeled [<region1>, <region2>, ...] in this image?

Table 9: A list of instructions for multiple-region caption.

1. Why?
2. What's the rationale for your decision
3. What led you to that conclusion?
4. What's the reasoning behind your opinion?
5. Can you explain the basis for your thinking?
6. What factors influenced your perspective?
7. How did you arrive at that perspective?
8. What evidence supports your viewpoint?
9. What's the logic behind your argument?
10. Can you provide some context for your opinion?
11. What's the basis for your assertion?
12. What experiences have shaped your perspective?
13. What assumptions underlie your reasoning?
14. What's the foundation of your assertion?
15. What's the source of your reasoning?
16. What's the motivation behind your decision?
17. What's the impetus for your belief?
18. What's the driving force behind your conclusion?
19. What's your reasoning?
20. What makes you say that?
21. What's the story behind that?
22. What's your thought process?
23. What's the deal with that?
24. What's the logic behind it?
25. What's the real deal here?
26. What's the reason behind it?
27. What's the rationale for your opinion?
28. What's the background to that?
29. What's the evidence that supports your view?
30. What's the explanation for that?

Table 10: A list of instructions for the second round chat in VCR.