

Shikra: Unleashing Multimodal LLM’s Referential Dialogue Magic

Keqin Chen^{12*}, Zhao Zhang^{1†}, Weili Zeng³, Richong Zhang², Feng Zhu¹, Rui Zhao¹⁴

¹SenseTime Research; ²SKLSDE, Beihang University

{³SEIEE, ⁴Qing Yuan Research Institute}, Shanghai Jiao Tong University

chenkq@act.buaa.edu.cn; zzhang@mail.nankai.edu.cn

Abstract

In human conversations, individuals can indicate relevant regions within a scene while addressing others. In turn, the other person can then respond by referring to specific regions if necessary. This natural referential ability in dialogue remains absent in current Multimodal Large Language Models (MLLMs). To fill this gap, this paper proposes an MLLM called Shikra, which can handle spatial coordinate inputs and outputs in natural language. Its architecture consists of a vision encoder, an alignment layer, and a LLM. It is designed to be straightforward and simple, *without* the need for extra vocabularies, position encoder, pre-/post-detection modules, or external plugin models. All inputs and outputs are in natural language form. Referential dialogue is a superset of various vision-language (VL) tasks. Shikra can naturally handle location-related tasks like REC and PointQA, as well as conventional VL tasks such as Image Captioning and VQA. Experimental results showcase Shikra’s promising performance. Furthermore, it enables numerous exciting applications, like providing mentioned objects’ coordinates in chains of thoughts and comparing user-pointed regions similarities. Our code and model are accessed at <https://github.com/shikras/shikra>.

1 Introduction

In recent months, Multimodal Large Language Models (MLLMs) have witness remarkable progress (Alayrac et al., 2022; Huang et al., 2023; Liu et al., 2023a; Zhu et al., 2023; Li et al., 2023a; Gao et al., 2023; Dai et al., 2023). They brings eyes to Large Language Models (LLMs), where users can talk about the input image. However, although these models can perceive image content, they cannot engage in dialogue with users regarding the

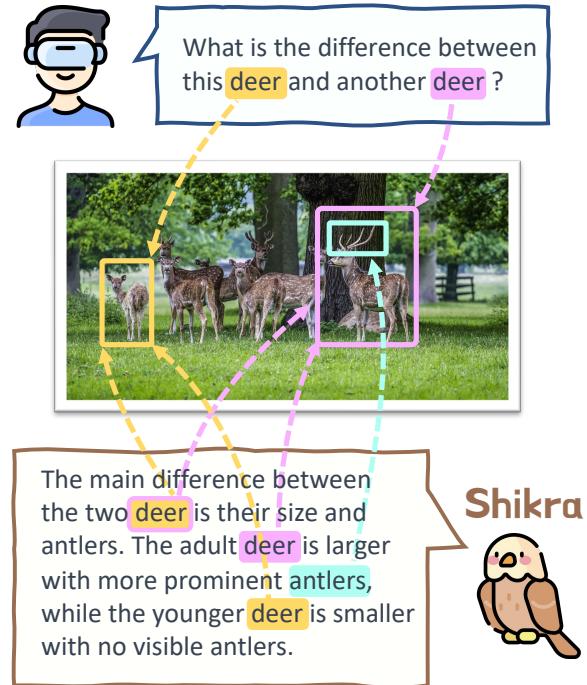


Figure 1: **Demo of Referential Dialogue (RD)**. Users can point to specific areas and ask questions. In turn, Shikra will indicate the specific regions when replying, if necessary. More interesting dialogues can be found in Figure 2 and Appendix C.

precise positions of the content. Users cannot indicate areas of interest in the image, and the models cannot provide the exact locations of the described content. Differently, as shown in Figure 1, in human daily communication, different regions or objects in the scene are often attended to, and people can speak and point to these regions for efficient information exchange. We refer to this interaction mode as Referential Dialogue (RD). If an MLLM excels in this skill, it will bring numerous exciting applications. For instance, applying it to Mixed Reality (XR) headsets like Apple Vision Pro, users can indicate anything to converse with the AI assistant. The AI assistant can display the prompt area in the field of view when necessary. It also assists

* Work done during internship at SenseTime Research

† Equal Contribution & Project Leader

visual robots in communicating with individuals by comprehending their specific reference positions. It facilitates online shopping by enabling users to inquire about items of interest in an image.

In this paper, we evolve MLLM to open the veil of referential dialogue. We create Shikra¹, a unified model capable of handling inputs and outputs of spatial coordinates. All coordinates, both input and output, are represented in natural language numerical form without introducing any extra vocabularies or position encoder. The Shikra architecture comprises a vision encoder, an alignment layer, and a LLM. We do **not** introduce any pre-/post-detection modules or external plug-in models, making Shikra unified and simple. We provide several real conversations with users in the Figure 2 and Appendix C, where users can use it to compare the differences of multiple regions, inquire about the meaning of the thumbnail, discuss specific objects, *etc.* Shikra can provide explanations when answering any question, not only verbally but also spatially.

Referential dialogue is a superset of many vision-language (VL) tasks. Shikra, skilled in RD, can naturally work on these tasks with promising performance, including Visual Question Answering (VQA), image captioning, and location-related tasks such as Referring Expression Comprehension (REC) and PointQA. We illustrate some of them in Figure 2. For more quantitative results, please refer to Section 6.3. Besides, this paper also addresses intriguing questions, such as how to represent position in an image (Section 6.2). Do previous MLLMs possess the capability to comprehend absolute positions? (Section 4). Can the reasoning process with location information assist in providing more accurate answers to questions? (Section 6.1). We hope that these analysis experiment can inspire future research on MLLMs.

The main contributions of this paper are:

- This paper introduces the task of Referential Dialogue (RD), which is an essential component of everyday human communication and possesses extensive practical applications.
- We present Shikra, a generalist MLLM, for RD. Shikra is simple and unified, **without** introducing extra vocabularies, pre-/post-detection module, or external plug-in models.

¹Shikra is a hunter’s companion, capable of understanding human language and gesture instructions, and locating and capturing prey in the wild.

- Shikra handles unseen settings effortlessly, creating diverse application scenarios. It also achieves promising performance on conventional visual language tasks such as REC, PointQA, VQA, and Image Captioning, without finetuning.

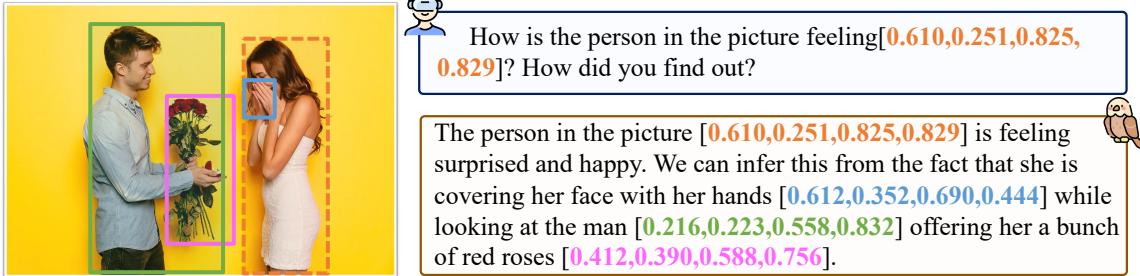
2 Related Works

2.1 Multimodal Large Language Model

Expanding the large language model to a multimodal version has garnered widespread attention. Flamingo (Alayrac et al., 2022) integrates visual adaption layers (like Perceiver) to an LLM, and trained on a large-scaled interleaved image-text dataset. OpenFlamingo (Awadalla et al., 2023) reimplements Flamingo and releases it to the community along with an M3C dataset. Subsequently, MM-GPT (Gong et al., 2023), and Otter (Li et al., 2023a) tune on carefully constructed instruction data for a more user-friendly interaction. Another genre is BLIP-2 (Li et al., 2023b), which align queried visual feature with text using multiple vision-language losses (model named Q-Former), and tunes a simple fully connection layer to feed the queried embedding to a frozen language model. Mini-GPT4 (Zhu et al., 2023), mPLUG-OWL (Ye et al., 2023), VPGTrans (Zhang et al., 2023a), and InstructBLIP (Dai et al., 2023) retain Q-Former, replace language model to a larger one, and then tuning on meticulously collected instruction data. Additionally, there are simpler and more direct methods: FROMAGe (Koh et al., 2023) and LLaVA (Liu et al., 2023a) directly feed visual features to the LLM using only a learnable fully connected layer. The closed source business model GPT-4 (OpenAI, 2023) also demonstrates astonishing image comprehension capabilities. Recently, interesting works have made remarkable progress by extending LLM to audio, *e.g.*, KOSMOS-1 (Huang et al., 2023), X-LLM (Chen et al., 2023), PandaGPT (Su et al., 2023) and control systems like PaLM-E (Driess et al., 2023) and EmbodiedGPT (Mu et al., 2023)

2.2 Vision-Language Positioning Tasks

Many vision-language tasks require localization representation. **Tasks with output boxes:** Referring Expression Comprehension (REC) (Kazemzadeh et al., 2014; Mao et al., 2016) aims to localize a target object in an image described by a referring expression.



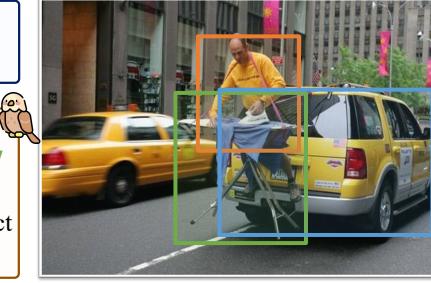
How is the person in the picture feeling [0.610,0.251,0.825,0.829]? How did you find out?



The person in the picture [0.610,0.251,0.825,0.829] is feeling surprised and happy. We can infer this from the fact that she is covering her face with her hands [0.612,0.352,0.690,0.444] while looking at the man [0.216,0.223,0.558,0.832] offering her a bunch of red roses [0.412,0.390,0.588,0.756].



What is unusual about this image? give coordinates [xmin,ymin,xmax,ymax] for the items you reference.



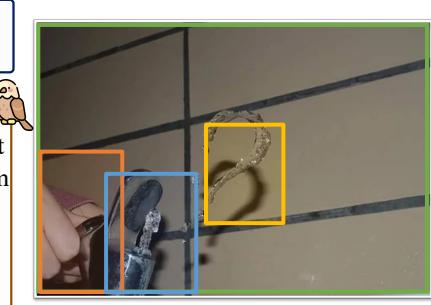
What is this man [0.171,0.330,0.650,0.779] scared of?



The man [0.171,0.330,0.650,0.779] is scared of the chicken [0.620,0.219,0.996,0.659] that is flying towards him.



Can you explain this meme? give coordinates [xmin,ymin, xmax,ymax] for the items you reference.



In this image, a person [0.002,0.490,0.208,0.832] is holding a water-spraying tool [0.180,0.546,0.408,0.830] and is pointing it at a tiled wall [0.002,0.168,0.998,0.830]. The water is dripping from the wall in the shape of the question mark [0.432,0.422,0.626,0.658]. This creates an interesting visual effect, as the question mark appears on the wall while the water is sprayed to resemble the question mark.

Figure 2: Referential Dialogues between real users and Shikra-7B. The dashed box on an image represents the area referred to by the user or jointly referred to by Shikra, while the solid box represents the area solely referred to by Shikra. More RD results and applications on conventional VL tasks can be found in Appendix C.

Described Object Detection (Xie et al., 2023) extends REC to more realistic scenarios where the object may not exist or there may be multiple objects. VQA Grounding aims to answer visual questions and associate the answers with specific visual regions or objects. **Tasks with input boxes:** Given an image and a location box, the task of Grounding Caption (GC) (Zhou et al., 2020) is to generate a description for this location by considering the surrounding environment. Compared to GC, Referring Expression Generation (REG) (Liu et al.,

2017) requires the generated description to indicate that it describes this region specifically, not others, making it necessary for the description to be discriminative. PointQA (Mani et al., 2020) requires a model answer for a visual question where the questioner queries a specific position in the picture. **Differently**, our model is not only compatible with the above tasks, but also can handle the input and output of position representation flexibly and simultaneously, bringing Referential Dialogue and extending new dimensions to positional tasks.

2.3 Position Representation

Inputting regions of interest into the model presents various approaches. Some methods (Bracha et al., 2023) directly concatenate cropped image patches with the original image as model input. There are also some methods (Lin et al., 2020, 2022) that use 0/1 mask or Gaussian map input with the original image to emphasize the area of user interest. Some methods (Tancik et al., 2020; Kirillov et al., 2023) first encode points and boxes to positional encodings then add them to intermediate features or learned queries. **Outputting** regions of interest is a highly focused technique, existing many positioning paradigms . Anchor-based methods utilize predefined sliding windows and proposal candidate regions for classification., e.g., Fast R-CNN (Girshick, 2015). Some one-stage methods remove anchors and directly regress four values for bounding box coordinates, e.g., FCOS (Tian et al., 2019). Some methods adopt one-to-one label assignment to evolve object detection into an end-to-end manner, e.g., DETR (Carion et al., 2020) and POTP (Wang et al., 2021). An interesting genre is Pix2seq (Chen et al., 2021), which formalizes the detection task as a sequence generation task. It desires the spatial position of the image in 1,000 bins and uses a 1,000-token vocabulary to represent it. For detection, Pix2seq performs classification on the coordinate vocabulary in an auto-regressive manner. Following Pix2seq, several methods, e.g., OFA (Wang et al., 2022b), Unified-IO (Lu et al., 2022), UniTab (Yang et al., 2022), GIT (Wang et al., 2022a), and VisionLLM (Wang et al., 2023b) introduce similar coordinate vocabulary alongside the language vocabulary for object detection and REC tasks. **Differently**, Shikra formulates position input/output as the most natural and flexible form of language and compare it with the extra coordinate vocabulary in Section 6.2.

3 Referential Dialogue

To better understand the interesting abilities of our model, we demonstrated real users’ communications in Figure 1 and Figure 2. As shown in the first demo of Figure 1, the user points to two deer, and inquires, “*What is the difference between this deer and another deer?*” When Shikra answered, she not only mention the differences but also output the coordinates of the differences. The subsequent examples in Figure 2 are alike. To our knowledge, there have been no unified models that can achieve

such functionality before. RD is a superset of numerous vision-language tasks. Shikra can perform most tasks like current MLLM, including VQA, Image Caption, and multimodal dialogue. Furthermore, it handles tasks that they cannot, like REC, REG, and PointQA. The model demonstrates proficiency in tasks not in the training set, such as identifying similarities between two indicated objects, or counting somethings, and providing their positions. We show more results in Appendix C. If you are interested in quantitative experiments, you can refer to Section 6 later.

4 Chessboard Test for Current MLLM

Can the current MLLM model understand absolute spatial positions? The current MLLMs cannot directly output coordinates; thus, in this section, we designed a chessboard test, which simplifies the object grounding into a part choice task. Specifically, we divide a image into a 2×2 chessboard. Next, we ask, “*<image> Which part is <expr> in if the picture is divided equally into four 2 by 2 parts? Choose from: (A) Top-left (B) Top-right (C) Bottom-left (D) Bottom-right.*”, where *<image>* and *<expr>* denote input image tokens and Class name. We construct test data from LVIS (Gupta et al., 2019), which is a perception detection with over 1000 entry-level object categories. We choose objects that are completely within a certain part (*i.e.*, ambiguous positions are not considered). In total, we select 600 images per part, resulting in 2,400 images across 945 categories. We employ LLaVA-13B (Liu et al., 2023a) for the chessboard test , but the results are unsatisfactory. We tried various instruction methods, and LLaVA should achieve an accuracy of 25.96%, which is comparable to random selection. This suggests that prior coarse-grained vision-language alignment pre-training may be inadequate for MLLMs to capture the exact spatial position of an image. We need to explore appropriate coordinate representations and finer-grained training data.

5 Breeding Shikra

This section introduces the birth of Shikra, encompassing its structure design, position representation, training data construction, and training strategies.

5.1 Architecture

We selected the pre-trained ViT-L/14 of CLIP as visual encoder and Vicuna-7/13B as our LLM. We

use one fully connected layer to map the ViT’s $16 \times 16 \times$ output embedding $\mathbf{V} \in \mathbb{R}^{16 \times 16 \times 1024}$ to $\mathbf{V}' \in \mathbb{R}^{256 \times D}$ for modal alignment and correct input dimension of LLM. D is 4,096 for Vicuna-7B and 5,120 for Vicuna-13B. Visual embedding can be inserted into anywhere of input sequence. During training, both the fully connected layer and the entire language model are involved. We do not introduce any vocabulary or special encoder for encoding position information. We have not introduced additional pre-/post-detectors for points or bounding boxes. The model using Vicuna-7B is called Shikra-7B, and the other, using Vicuna-13B, is named Shikra-13B.

5.2 Numerical representation of position

We represent the position using numerical values in Natural Language in a highly intuitive manner. We use $[x_{\min}, y_{\min}, x_{\max}, y_{\max}]$ to denote the bounding box and $[x_{\text{center}}, y_{\text{center}}]$ to denote region center point. x and y is normalized according to the size of the image. We default to keeping 3 decimal places for each number. These coordinates can appear anywhere in the input and output sequence of the model. For example, User Question: “*How many other clothes in the <image> are of the same color as the jacket [0.268, 0.372]?*”. Shikra reply: “*The jacket [0.268, 0.372] is green. We can find a T-shirt [0.653, 0.532] and cropped pants [0.569, 0.101] a with same green color. So the answer is two.*” The square brackets that record coordinates naturally appear in sentences and can serve as any sentence component. Like regular text, tokenizing without discrimination.

5.3 Instruction data construction

We utilize two types of data to train Shikra: the reorganized public datasets, and the high-quality RD data built from Flickr30K Entities (Plummer et al., 2015) using GPT-4 (OpenAI, 2023).

5.3.1 Reorganization of public data

We collection training data from public VQA, Image Captioning datset, and several datasets already containing positional annotation, such as RefCOCO (Kazemzadeh et al., 2014) for REC/REG, visual gemone (Krishna et al., 2017) for grounding caption, Visual-7W (Mani et al., 2020) for PointQA. We also define new task forms, such as Spotting Captioning, which requires the model to describe the image and spots the mentioned objects or regions using points or boxes. We use Flickr30K En-

tities for this task. All the data used and corresponding tasks can be found in Appendix A. Note that all the data used were included in the reported model results, unless stated otherwise for specific comparative experiments. Additionally, it should be mentioned that we have **excluded** images present in the test and validation data from the training data to prevent potential data leakage, despite their distinction in terms of image-text pairs.

5.3.2 Generated data

The existing publicly available data is not sufficient to train an MLLM skilled in RD, as they lack CoT data with positional annotations, natural communication data with positional annotations, etc. We resort to GPT-4 to obtain high-quality RD annotations from Flickr30K Entities. Flickr30K Entities has five descriptions for each image. These mentioned objects appearing in the image will be labeled using bounding box. Although the API of GPT-4 temporarily **cannot** see images, we explained the format of the bounding boxes to GPT-4 and asked it to understand the image through these five sentences and boxes. Next, we require GPT-4 to design Q&A pairs. When designing problems, these questions must be able to determine answers from known information. In this way, we generated 5,922 QA pairs, where coordinate information may appear in both questions and answers. The dataset will continue expanding in the future. You can refer to it as Shikra-RD.

5.3.3 Task prompts

We construct variable task templates for different tasks. For instance, for the spottingS caption task, we can use “*Can you provide a description of the image <image> and include the coordinates [x0,y0,x1,y1] for each mentioned object?*” where <image> represents the visual tokens. For PointQA, we can use “*Referring to point <objs> in image <image>, give a direct answer to ’<question>’*” where <objs> denotes the coordinates of the region and <question> represents the question from the source dataset. For REC, “*In <image>, I need the bounding box coordinates of <expr>.*” where <expr> is the expression. More templates for different tasks can be found in the Appendix.

It should be noted that we cannot use an invariant task template for a specific type of task. In this case, the model cannot flexibly accept user instructions. To solve this problem, we first describe the purpose of specific tasks, write a sample template, and then

Table 1: **Comparing different forms of CoTs.** We train three toy models of Shikra-7B (without using additional datasets) on the CLEVR dataset. Q, A, C, and C^{Point} denote the Question, final Answer, Chain of thoughts, and Chain of thoughts with Pointing.

$Q \rightarrow A$	$Q \rightarrow CA$	$Q \rightarrow C^{\text{Point}} A$
88.07	80.68	93.97

have GPT-4 rewrite it in rich language, expanding it into hundreds of variations to convey the same meaning. During training, we can randomly choose from them. We provide details on some generated task templates in the Appendix B.

5.4 Tuning details

Shikra is trained in two stages. In the first stage, we train it on the reorganized VL dataset (Section 5.3.1) for 100,000 steps (around 1.5 epoch); In the second stage, we raise the sampling ratio to 50% on LLaVA-Instruct-150K (Liu et al., 2023a) and our generated RD data (Section 5.3.2). In both stages, we freeze the visual encoder and tune all parameters in LLM. We adopt AdamW (Loshchilov and Hutter, 2019) as the optimizer and cosine annealing scheduler (Loshchilov and Hutter, 2017) as learning rate scheduler with an initial learning rate of 2e-5 and global batch size of 64. All training runs on 8 NVIDIA A100 GPUs. It takes around 100h for stage one training and 20h for stage two.

6 Experiment and Analysis

6.1 Grounding CoT or verbal CoT?

The process of providing reasoning before giving an answer is called Chain of the thoughts (CoT), which provides good explanatory during model judgments. However, CoT often suffer from hallucinations (Zhang et al., 2023b), which often do not improve the performance of the final answer. Current MLLMs are also suffer from serious visual hallucination (Li et al., 2023c). In this section, we investigate whether CoT with position annotations can reduce hallucinations and improve model performance. In this paper, we refer to this type of CoT as Grounding CoT (GCoT). We train our Shikra-7B (without pre-training) on CLEVR (Johnson et al., 2017) in three settings: 1) Only use Question and Answer ($Q \rightarrow A$); 2) Use Question, CoT, and answer ($Q \rightarrow CA$); 3) Use GCoT with Center Point annotation and answer ($Q \rightarrow C^{\text{Point}} A$). We record they

Table 2: **Comparing different position representations.** We implement Shikra-7B in two different representation forms and train two toy models solely on RefCOCO, RefCOCO+/g, and Visual Genome for controllable comparison. Vocab. means to use extra vocabularies to represent coordinates, like (Chen et al., 2021; Wang et al., 2022b), and Numerical means to directly use numerals in natural language to express coordinates.

Dataset	Split	Vocab.	Numerical
RefCOCO	val	81.03	81.47
	test-A	86.94	87.40
	test-B	70.91	73.25
RefCOCO+	val	72.32	74.30
	test-A	81.78	83.29
	test-B	59.95	63.08
RefCOCOg	val-u	72.81	75.69
	test-u	73.78	75.52

performance in Table 1. Using only CoT to train the model ($Q \rightarrow CA$) and requiring a reasoning process before the final answer decreases performance compared to direct answering setting ($Q \rightarrow A$). In the $Q \rightarrow C^{\text{Point}} A$ setting, we ask the model to provide CoT along with center points $[x_{\text{center}}, y_{\text{center}}]$ for each mentioned object. Performance improved by 13 points compared to $Q \rightarrow CA$ and 5.9 points compared to $Q \rightarrow A$, indicating that training with positional annotations suppresses visual hallucination. This is a preliminary attempt at GCoT, and it is a promising direction worth exploring.

6.2 Location tokens or just numbers?

For detect object in autoregressive model, several methods (Chen et al., 2021; Wang et al., 2022b) introduce extra vocabularies (e.g., $\langle \text{bin_}0 \rangle, \dots, \langle \text{bin_}1000 \rangle$) to represent coordinates for object detection in spatially discretized images, as described in Section 2.3. In contrast, Shikra represents coordinates naturally and intuitively, using numbers directly. Which form is better? We train two toy Shikra using two different representations with REC data, they performance is recorded in Table 2, where using numbers directly achieves better results. Aside from performance, our simple-designed coordinate numerical representation makes the model more elegant without modifying vocabularies for localization tasks. Users can freely control the precision of numerical representation (number of digits after the decimal separator) without retraining vocabularies. How-

Table 3: **Results on standard REC task.** Generalist VL model Generalist VL models can directly perform various vision-language tasks, including image captioning, VQA, REC, etc. Specialist models are those specifically designed for localization tasks (e.g., UNINEXT, Yan et al., 2023 and G-DINO, Liu et al., 2023b), or generalist pretraining models that have undergone multitask localization finetuning (e.g., Yang et al., 2022) or single-task finetuning (e.g., Wang et al., 2022b). We select the three current best performing models (Liu et al., 2023b; Yan et al., 2023; Wang et al., 2023a) as baselines. OFA-L* (Wang et al., 2022b) refers to the OFA-Large checkpoint without finetuning. GRIT refexp is the ablation split (Lu et al., 2022).

Model type	Model	RefCOCO			RefCOCO+			RefCOCOg		GRIT refexp
		val	test-A	test-B	val	test-A	test-B	val-u	test-u	
Generalist VL SOTAs (w/o finetuning)	GPV-2	-	-	-	-	-	-	-	-	51.50
	OFA-L*	79.96	83.67	76.39	68.29	76.00	61.75	67.57	67.58	61.70
	Unified-IO	-	-	-	-	-	-	-	-	78.60
	OFASys	-	80.10	-	-	-	-	-	-	-
	VisionLLM-H	-	86.70	-	-	-	-	-	-	-
	Shikra-7B	87.01	90.61	80.24	81.60	87.36	72.12	82.27	82.19	69.34
Specialist SOTAs (Specialist/Finetuned)	Shikra-13B	87.83	91.11	81.81	82.89	87.79	74.41	82.64	83.16	69.03
	G-DINO-L	90.56	93.19	88.24	82.75	88.95	75.92	86.13	87.02	-
	UNINEXT-H	92.64	94.33	91.46	85.24	89.63	79.79	88.73	89.37	-
	ONE-PEACE	92.58	94.18	89.26	88.77	92.21	83.23	89.22	89.27	-

Table 4: **Comparing pointQA capabilities on the Visual-7W** (Zhu et al., 2016). Visual-7W features a ‘which box’ setting, requiring the model to select one matching box from four options based on the given description. Accuracy (%) is used for evaluation.

Zhu et al.	Hu et al.	Lu et al.	Lu et al.*	Shikra
56.10	72.53	82.75	83.35	85.33

ever, it also has drawbacks. Compared to using extra vocabularies, numerical representation requires more tokens to represent coordinates, leading to increased computational costs when predicting dense objects. In this paper, we still prefer numerical representation, but future research can choose the appropriate method based on their pros and cons.

6.3 Quantitative results on conventional tasks

Our Shikra excels in Referential Dialogue, facilitating seamless integration into a wide range of vision-language (VL) tasks, particularly those related to positioning. Here, we present the quantitative results for these tasks.

To demonstrate the positioning capability of our model, we examine the REC task, in which models are asked to ground the object described with an expression. As shown in Table 3, we compare our method with generalist VL models that perform multiple tasks without finetuning. We also compare our method with Specialist SOTAs, including localization specialist models and gen-

Table 5: **Comparing pointQA capabilities on the LookTwice-QA** (Mani et al., 2020), where the models are asked to answer question based on the input point/box. Pronoun, Superclass (Super cls.), and Class indicate different levels of referential clarity in the question, e.g., “How many of these [∅/fruits/apples] <obj>?” We use Shikra-13B and Accuracy (%) for evaluation.

Type	Point		Box		
	Model	Mani et al.	Shikra	Model	Shikra
Pronoun		56.5	70.0	60.2	70.3
Super cls.		59.1	70.2	59.8	71.4
Class		62.8	71.8	61.4	72.3

alist/foundation models that perform specific finetunes on localization-related tasks. In this setting, we instruct Shikra to provide the coordinates of the objects referred to by the expression. For an example, we use “I’d like to know the exact coordinates of <expr> in the photo <image>.”, where <expr> represents the expression and <image> represents the input image. More instructions can be found in Appendix 9. The experimental results demonstrate that Shikra achieves promising performance compared to other generalist models.

Correspondingly, to quantitatively evaluate our model’s understanding of position inputs, we evaluated our model on two types PointQA datasets, LookTwice-QA of (Mani et al., 2020) and Visual7W (PointQA Setting) of (Zhu et al., 2016).

Table 6: **Comparing generalist models on VQA and Image Captioning.** For VQA, we evaluate SOTA generalist models and our Shikra-13B on VQAv2 (Antol et al., 2015) and OK-VQA (Marino et al., 2019) following the normalization rules. Here, we also provide VQAv2^{val} (83.3) and OK-VQA (53.8) results on LVLM-eHub toolbox (Xu et al., 2023) for easy comparison. For Image Captioning, we evaluate them on COCO (Chen et al., 2015) and Flickr30k (Plummer et al., 2015) in CIDEr. We call Flamingo (Alayrac et al., 2022) FM for short.

Datasets		Shikra	FM-80B	FM-9B	Kosmos-1	BLIP-2	Unified-IO	VPGTrans	VisionLLM
VQA	VQAv2 ^{val}	75.33	-	-	-	65.2	-	65.2	-
	VQAv2 ^{dev}	77.36	56.3	51.8	51.0	65.0	77.9	-	-
	VQAv2 ^{std}	77.51	-	-	-	-	-	-	-
	OK-VQA	47.16	50.6	44.7	-	45.9	54.0	45.0	-
Caption	Flickr30k	73.9	67.2	61.5	67.1	-	-	-	-
	COCO	117.5	84.3	79.4	84.7	-	122.3	-	114.2

Table 7: **Object hallucination benchmark using POPE evaluation pipeline** (Li et al., 2023c). Accuracy denotes the accuracy of predictions. Precision signifies the true positive samples among the predicted positives. Recall indicates the correct identification of all true positive samples. “Yes” represents the probability of the model outputting a positive answer. Except for Shikra-7B, the other results are obtained from Li et al., 2023c.

Datasets	Metrics	Shikra	InstructBLIP	MiniGPT-4	LLaVA	MM-GPT	mPLUG-Owl
Random	Accuracy (↑)	86.90	88.57	79.67	50.37	50.10	53.97
	Precision (↑)	94.40	84.09	78.24	50.19	50.05	52.07
	Recall (↑)	79.27	95.13	82.20	99.13	100.00	99.60
	F1-Score (↑)	86.19	89.27	80.17	66.64	66.71	68.39
	Yes	43.26	56.57	52.53	98.77	99.90	95.63
Popular	Accuracy (↑)	83.97	82.77	69.73	49.87	50.00	50.90
	Precision (↑)	87.55	76.27	65.86	49.93	50.00	50.46
	Recall (↑)	79.20	95.13	81.93	99.27	100.00	99.40
	F1-Score (↑)	83.16	84.66	73.02	66.44	66.67	66.94
	Yes	45.23	62.37	62.20	99.40	100.00	98.57
Adversarial	Accuracy (↑)	83.10	72.10	65.17	49.70	50.00	50.67
	Precision (↑)	85.60	65.13	61.19	49.85	50.00	50.34
	Recall (↑)	79.60	95.13	82.93	99.07	100.00	99.33
	F1-Score (↑)	82.49	77.32	70.42	66.32	66.67	66.82
	Yes	46.50	73.03	67.77	99.37	100.00	98.67

LookTwice-QA asks models to answer questions about the region specified by the user, either by center point or box, with the distinction that these questions necessitate comprehending the user-designated area first, and then observing the entire image to answer. For instance, “*How many of these [Pronoun/Superclass/Class] <obj>?*”, where <obj> denotes the coordinates of input point or box and [Pronoun/Superclass/Class] represents language instructions with different clarity levels (e.g., [\emptyset /fruits/apples]). Visual7W also provides a setting for point QA, where models are given a question and four box options, and should choose one as the answer.’ Our Shikra achieves the SOTA per-

formance in all these settings.

Additionally, we assess our model on conventional VL tasks in Table 6, such as VQA and Image Captioning, which do not necessitate coordinates in their input or output. The experimental results show that we achieved promising results on most datasets. We also evaluated the performance of our method in POPE evalution pipeline (Li et al., 2023c), and the results are recorded in Table 7. Our method has achieved results comparable to InstrutBLIP(Dai et al., 2023) and far surpasses recent popular MLLMs. It’s worth noting that these task configurations are just some subsets of Referential Dialogue. We hope readers can appreciate the

more intriguing capabilities of Shikra in Figure 2 and Appendix C.

7 Limitations

Shikra only supports English and is not user-friendly for non-English speakers. Making Shikra multilingual in the future is valuable. Shikra is unsuitable for dense object detection and segmentation tasks. Exploring improved coordinate representations for these tasks is also interesting. Shikra, like most LLMs, may produce harmful and counterfactual responses.

8 Conclusion

Our study unveiled the critical gap in MLLMs' ability to understand and engage in referential dialogue, an integral aspect of human communication. To address this, we introduced Shikra, a unified, straightforward model designed to comprehend and output spatial coordinates in natural language. Our approach does not necessitate extra vocabularies, position encoders, or external plug-ins, preserving the model's simplicity. It was proved that Shikra performs notably well on a variety of conventional vision-language tasks, while offering swathes of exciting applications such as aiding AI assistants in Mixed Reality headsets or facilitating precise communication in online shopping scenery.

References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Anas Awadalla, Irena Gao, Joshua Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Open-flamingo.
- Lior Bracha, Eitan Shaar, Aviv Shamsian, Ethan Fetaya, and Gal Chechik. 2023. Disclip: Open-vocabulary referring expression generation. *arXiv preprint arXiv:2305.19108*.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer.
- Feilong Chen, Minglun Han, Haozhi Zhao, Qingyang Zhang, Jing Shi, Shuang Xu, and Bo Xu. 2023. X-LLM: Bootstrapping advanced large language models by treating multi-modalities as foreign languages. *arXiv preprint arXiv:2305.04160*.
- Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. 2021. Pix2seq: A language modeling framework for object detection. *arXiv preprint arXiv:2109.10852*.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructclip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. 2023. PaLM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*.
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aoju Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. 2023. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*.
- Ross Girshick. 2015. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448.
- Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv preprint arXiv:2305.04790*.
- Agrim Gupta, Piotr Dollar, and Ross Girshick. 2019. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364.
- Ronghang Hu, Marcus Rohrbach, Jacob Andreas, Trevor Darrell, and Kate Saenko. 2017. Modeling relationships in referential expressions with compositional modular networks. In *CVPR*, pages 1115–1124.

- Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. 2023. Language is not all you need: Aligning perception with language models. *arXiv preprint arXiv:2302.14045*.
- Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. 2017. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, pages 2901–2910.
- Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. 2014. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, pages 787–798.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. *arXiv preprint arXiv:2304.02643*.
- Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. 2023. Grounding language models to images for multimodal generation. *arXiv preprint arXiv:2301.13823*.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual Genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123:32–73.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023c. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
- Zheng Lin, Zhao Zhang, Lin-Zhuo Chen, Ming-Ming Cheng, and Shao-Ping Lu. 2020. Interactive image segmentation with first click attention. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13339–13348.
- Zheng Lin, Zhao Zhang, Ling-Hao Han, and Shao-Ping Lu. 2022. Multi-mode interactive image segmentation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 905–914.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Jingyu Liu, Liang Wang, and Ming-Hsuan Yang. 2017. Referring expression generation and comprehension via attributes. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4856–4864.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023b. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*.
- Ilya Loshchilov and Frank Hutter. 2017. **SGDR: stochastic gradient descent with warm restarts**. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Ilya Loshchilov and Frank Hutter. 2019. **Decoupled weight decay regularization**. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. 2022. Unified-dio: A unified model for vision, language, and multi-modal tasks. *arXiv preprint arXiv:2206.08916*.
- Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 2020. 12-in-1: Multi-task vision and language representation learning. In *CVPR*, pages 10437–10446.
- Arjun Mani, Nobline Yoo, Will Hinthon, and Olga Russakovsky. 2020. Point and ask: Incorporating pointing into visual question answering. *arXiv preprint arXiv:2011.13681*.
- Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. 2016. Generation and comprehension of unambiguous object descriptions. In *CVPR*, pages 11–20.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204.
- Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhui Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. 2023. Embodiedgpt: Vision-language pre-training via embodied chain of thought. *arXiv preprint arXiv:2305.15021*.
- OpenAI. 2023. **Gpt-4 technical report**.
- Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649.

- Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. 2023. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*.
- Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. 2020. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems*, 33:7537–7547.
- Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. 2019. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9627–9636.
- Jianfeng Wang, Lin Song, Zeming Li, Hongbin Sun, Jian Sun, and Nanning Zheng. 2021. End-to-end object detection with fully convolutional network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15849–15858.
- Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. 2022a. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*.
- Peng Wang, Shijie Wang, Junyang Lin, Shuai Bai, Xiaohuan Zhou, Jingren Zhou, Xinggang Wang, and Chang Zhou. 2023a. ONE-PEACE: Exploring one general representation model toward unlimited modalities. *arXiv preprint arXiv:2305.11172*.
- Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022b. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR.
- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. 2023b. Visionilm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*.
- Chi Xie, Zhao Zhang, Yixuan Wu, Feng Zhu, Rui Zhao, and Shuang Liang. 2023. Exposing the troublemakers in described object detection. *arXiv preprint*.
- Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. 2023. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *arXiv preprint arXiv:2306.09265*.
- Bin Yan, Yi Jiang, Jiannan Wu, Dong Wang, Ping Luo, Zehuan Yuan, and Huchuan Lu. 2023. Universal instance perception as object discovery and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15325–15336.
- Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Faisal Ahmed, Zicheng Liu, Yumao Lu, and Lijuan Wang. 2022. Unitab: Unifying text and box outputs for grounded vision-language modeling. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVI*, pages 521–539. Springer.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*.
- Ao Zhang, Hao Fei, Yuan Yao, Wei Ji, Li Li, Zhiyuan Liu, and Tat-Seng Chua. 2023a. Transfer visual prompt generator across llms. *arXiv preprint arXiv:2305.01278*.
- Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023b. Multi-modal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*.
- Yuanen Zhou, Meng Wang, Daqing Liu, Zhenzhen Hu, and Hanwang Zhang. 2020. More grounded image captioning by distilling image-text matching model. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4777–4786.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.
- Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4995–5004.

A Details of All Training Data

We listed all training data in Table 8. The asterisk indicates that this data is only used in the second training stage. We removed the images from the training set that are the same as those in the testing or validation set to prevent potential data leakage.

B Examples of Task Prompts

We list some task prompts used by Shikra during training in Table 9. For every task listed, there are hundreds. These prompts are generated by GPT-4 with carefully designed instructions. We randomly

Table 8: **All training data used by Shikra.** The asterisk indicates that this data is only used in the second stage.

Task	Dataset
Captioning	LLaVA-Pretraining
Spotting Cap.	Flickr30K Entities
Grounding Cap.	Visual Genome
REG	RefCOCO, RefCOCO+, RefCOCOg
REC	RefCOCO, RefCOCO+, RefCOCOg, Visual Genome
VQA	VQAv2
PointQA	PointQA-Local/Twice, Visual-7W ('which box' subset)
Dialogue	LLaVA-Instruct-150K*
RD	VCR, Shikra-RD (Generated data from Flickr30K Entities)*

Table 9: **Examples of task templates used by Shikra on different types of training data.** The explanation of placeholders in the template is as follows: ‘<image>’ represents the input image; ‘<objs>’ refers to the center points or bounding box of a user-specified location; ‘<question>’ denotes the question in the VQA dataset; ‘<expr>’ represents the expression in the REC task. During inference, there is no need to be confined to these forms. Users can describe their needs in natural language, creating more diverse and engaging task formats.

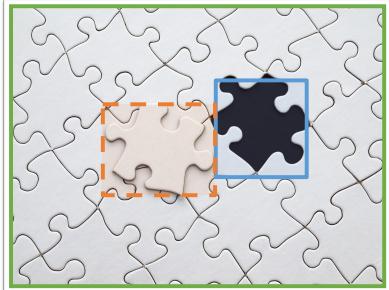
Task	Three randomly chosen examples from hundreds.
Captioning	Describe this image <image> as simply as possible. What is the content of the image <image>? Please answer in short sentences. Summarize the content of the photo <image>.
Spotting Cap.	Can you provide a description of the image <image> and include the coordinates [x0,y0,x1,y1] for each mentioned object? Please explain what's happening in the photo <image> and give coordinates [xmin,ymin,xmax,ymax] for the items you reference. How would you describe the contents of the image <image>? Please provide the positions of mentioned objects in square brackets.
Grounding Cap.	Can you give me a description of the region <objs> in image <image>? Describe what's happening within the coordinates <objs> of the given image <image>. What does the area <obj> within the given visual <image> contain?
REG	For the given image <image>, can you provide a unique description of the area <objs>? In the photo <image>, how would you describe the selected area <obj> uniquely? Can you provide a description for the region <objs> in the image <image> such that it sets it apart from others?
Q → A	I want to know the answer to ‘<question>’ Refer to the image <image> and give a clear response. Answer this question directly after referring to the image <image>: <question> Examine the image <image> and provide a brief answer for ‘<question>’
Q → CA	Having a look at image <image>, can you tell me the answer to my question ‘<question>’ and the logic leading to it? Please answer the following question ‘<question>’ based on the image <image>, and describe your thought process Upon analyzing the image <image>, please find the answer to my question ‘<question>’ and provide a detailed explanation.
Q → C ^{Point} A	Analyze the image <image> and answer ‘<question>’ Include your reasoning process and mark center points of related objects as [cx, cy]. Based on <image>, please respond to ‘<question>’ Include your thought process and note involved objects using [cx, cy] for their center points. While observing image <image>, kindly answer ‘<question>’ Elaborate on your reasoning process and tag any object center points involved [x,y].
Q → C ^{Box} A	<question> Please offer your reasoning process, and provide bounding boxes of mentioned objects within square brackets. Here is the picture <image> Please explain your reasoning and provide bounding boxes, denoted by square brackets, for the objects mentioned in the picture <image>. <question> Consider the image <image>, and then provide a well-reasoned answer to the question ‘<question>’ Don’t forget to mark relevant object locations using [x0,y0,x1,y1].
REC	In the given <image>, could you find and tell me the coordinates of <expr>? I need the coordinates of <expr> in <image>, can you please assist me with that? Locate <expr> in <image> and provide its coordinates, please.

selected three prompts for readers’ better understanding. Note that during inference, there is no need to confine oneself to these forms. Users can express their needs in natural language, creating diverse and engaging task formats.

C More Conversations with Shikra

We provide additional dialogue records of Shikra-7B in this section. For instance, we showcase RD results in Figure 3, VQA (Q→C^{Box} A) in Figure

4, and Spotting Captioning in Figure 6. We also include examples of traditional VL task forms, like OCR in Figure 5, REC in Figure 8, REG in Figure 7, and PointQA in Figure 9. Furthermore, Figure 9 and Figure 10 demonstrates that our input and output can handle points and boxes, just tell Shikra what to do.



Referential Dialogue

How can I complete this puzzle[**0.251,0.385,0.548,0.630**]?

To complete the puzzle[**0.002,0.126,0.998,0.872**], you need to place the missing piece[**0.251,0.385,0.548,0.630**] into the empty space[**0.548,0.322,0.790,0.580**].

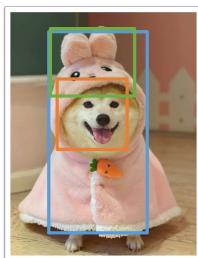


Referential Dialogue

Which one is the woman's reflection in the mirror?
[**0.594,0.270,0.797,0.756;0.369,0.276,0.572,0.755**]

The woman [**0.594,0.270,0.797,0.756**] is standing in front of the glass doors [**0.290,0.170,0.578,0.830**], which create a reflection [**0.369,0.276,0.572,0.755**] in the mirror.

Figure 3: **Referential Dialogue using Shikra-7B.** The dashed box on an image represents the area referred to by the user or jointly referred to by Shikra, while the solid box represents the area solely referred to by Shikra.



$\mathbf{Q} \rightarrow \mathcal{C}^{\text{Box}} \mathbf{A}$

Is this a rabbit?

The animal is a dog [**0.316,0.278,0.600,0.562**], not a rabbit. Although it is wearing a pink hooded blanket [**0.282,0.088,0.680,0.898**] and a bunny-like hat [**0.292,0.080,0.630,0.348**], it is a dog. Therefore, the answer is no.



What is this animal?

$\mathbf{Q} \rightarrow \mathcal{C}^{\text{Box}} \mathbf{A}$



The animal in the image is a giraffe [**0.350,0.062,0.600,0.916**]. It has a very distinctive appearance, with a mix of human clothing and characteristics.

Figure 4: **$\mathbf{Q} \rightarrow \mathcal{C}^{\text{Box}} \mathbf{A}$ using Shikra-7B.** It asks models generate grounded explanation for the answer.



What is written on this board?[**0.656,0.475,0.888,0.584**]



The word "Change" is written on the board [**0.656,0.475,0.888,0.584**].

Figure 5: **OCR using Shikra-7B.** We do not have explicit OCR datasets in Shikra training.

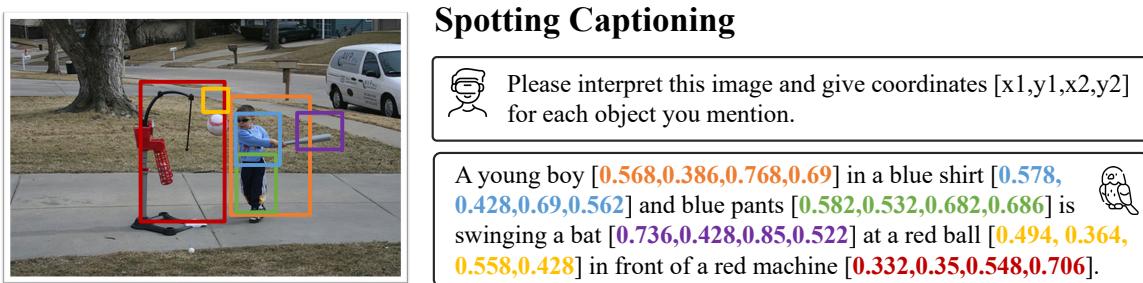
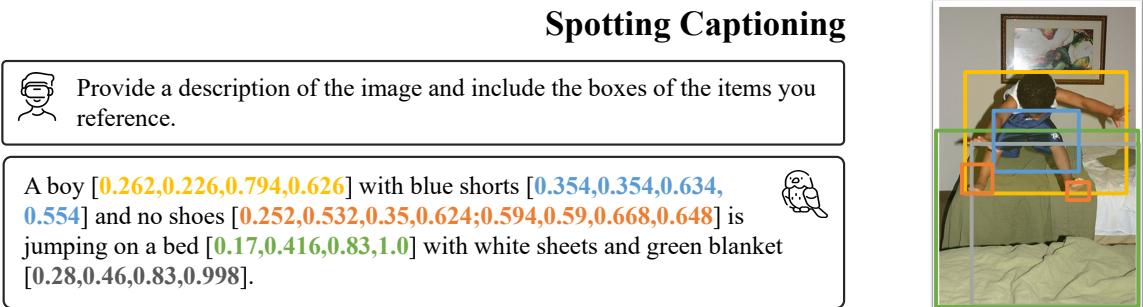
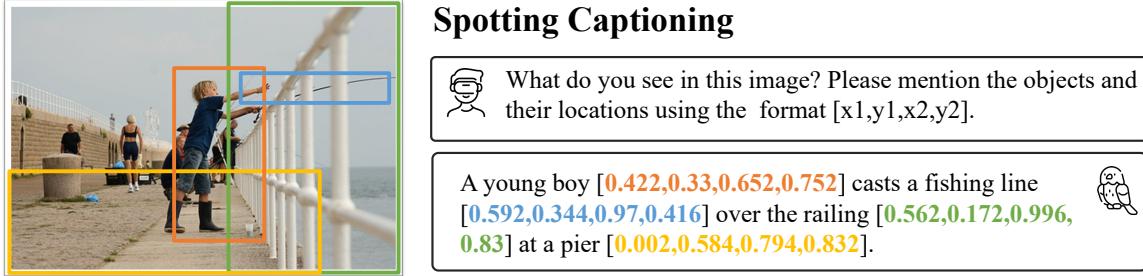


Figure 6: **Spotting Captioning using Shikra-7B.** The task requires the model to describe the image and spots the mentioned objects or regions using points or boxes.

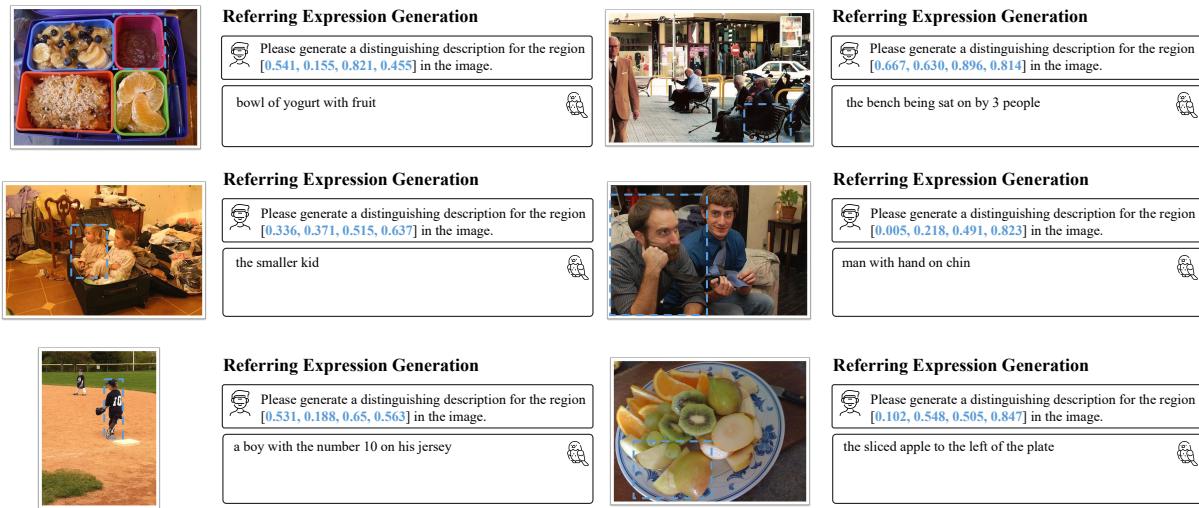


Figure 7: **Referring Expression Generation (REG) using Shikra-7B.** The purpose of REG is to generate a unique description for a specified location.

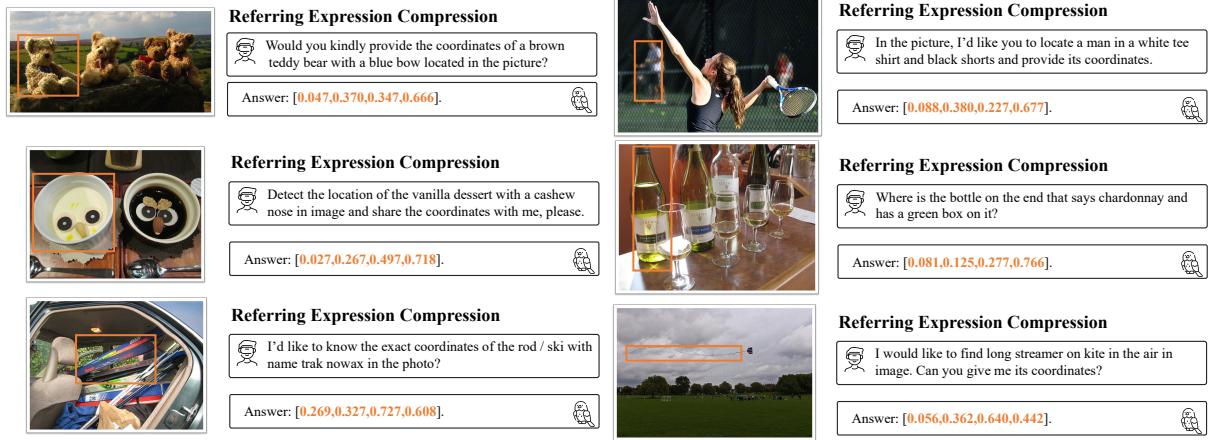


Figure 8: **Referring Expression Comprehension (REC) using Shikra-7B.** The task aims to localize a target object in an image described by a referring expression.

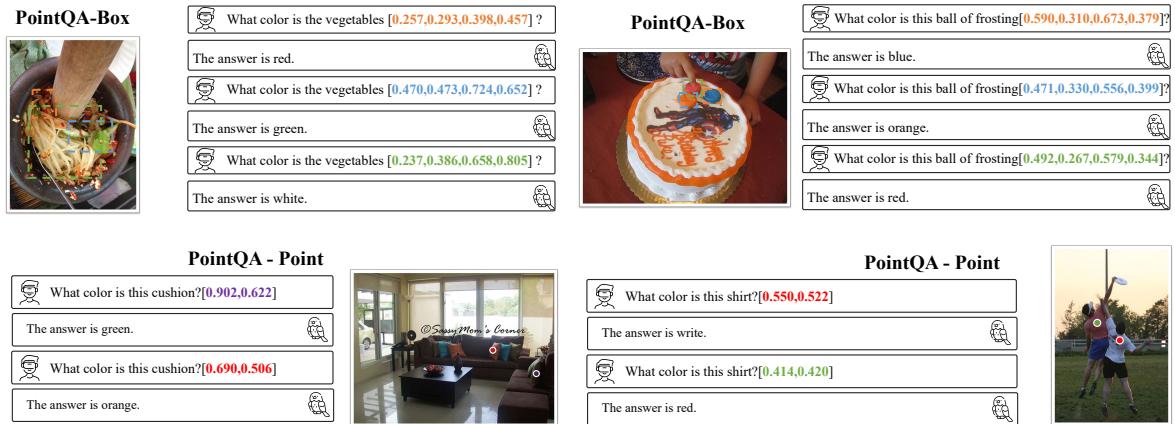


Figure 9: **PointQA using Shikra-7B.** The task asks models to answer questions about the region specified by the user, either by center point or box.

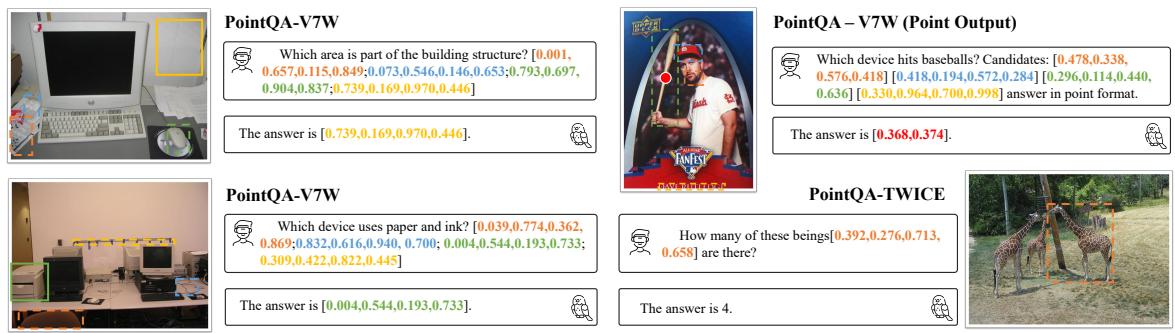


Figure 10: **PointQA-V7W using Shikra-7B.** PointQA-V7W provides a setting for point QA, where models are given a question and four box options, and should choose one as the answer.