

NExT-Chat: An LMM for Chat, Detection and Segmentation

Ao Zhang
National University of Singapore
aozhang@u.nus.edu

Chen-Wei Xie
Alibaba Group
eniacycw@alibaba-inc.com

Wei Ji
National University of Singapore
jiwei@nus.edu.sg

Liming Zhao
Alibaba Group
lingchen.zlm@alibaba-inc.com

Yun Zheng
Alibaba Group
zhengyun.zy@alibaba-inc.com

Tat-Seng Chua
National University of Singapore
dcscts@nus.edu.sg

<https://next-chatv.github.io>

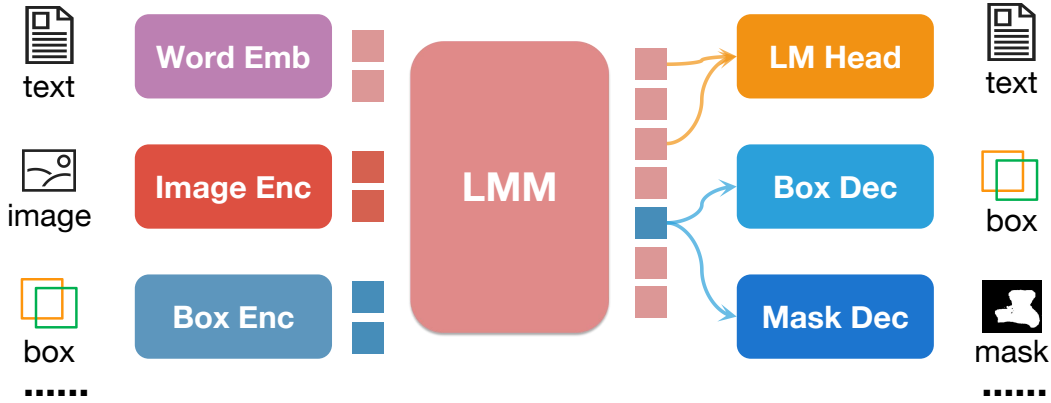


Figure 1. By using the embedding based location modeling method, the LMM can take bounding boxes as input and output text, bounding boxes and masks in the multimodal conversation.

Abstract

The development of large language models (LLMs) has greatly advanced the field of multimodal understanding, leading to the emergence of large multimodal models (LMMs). In order to enhance the level of visual comprehension, recent studies have equipped LMMs with region-level understanding capabilities by representing object bounding box coordinates as a series of text sequences (pixel2seq). In this paper, we introduce a novel paradigm for object location modeling called pixel2emb method, where we ask the LMM to output the location embeddings and then decode them with different decoders. This paradigm allows us to use different location formats (such as bounding boxes and masks) in multimodal conversations. Furthermore, this kind of embedding based location modeling enables the utilization of existing practices in localization tasks, such as

detection and segmentation. In scenarios with limited resources, our pixel2emb demonstrates superior performance compared to existing state-of-the-art (SOTA) approaches in both the location input task (72.4 vs. 71.2) and output task (79.4 vs. 71.4) under fair comparison. Leveraging the proposed pixel2emb method, we train an LMM named NExT-Chat and demonstrate its capability of handling multiple tasks like visual grounding, region caption, and grounded reasoning. We will release all the codes and models associated with this research.

1. Introduction

Recently, large language models (LLMs) have shown spreading influence in different areas, among which large multimodal models (LMMs) is one of the most attractive

area. Researchers try to equip LLMs with visual perception modules resulting in LMMs [6, 12, 27, 29] that can describe the visual content and answer visual questions. However, these LMMs are limited to holistic image understanding without the ability to conduct region-level reasoning, for example, locating the referred objects in the conversation.

To enable region-level understanding, current solutions [4, 17, 22] utilize the pixel2seq [5] paradigm where the object coordinates are converted to LLM understandable text tokens (e.g., $[x_1, y_1, x_2, y_2]$). Consequently, LMMs can output object coordinates as part of a normal next token prediction problem. However, the pixel2seq paradigm is limited to discrete coordinate outputs and struggles to provide other fine-grained formats, such as segmentation masks.

To address these limitations, we propose the pixel2emb paradigm, which can accommodate different location formats. Specifically, we introduce two new tokens, `<trigger>` and `<loc>`, where the `<trigger>` serve as a trigger for localization and `<loc>` act as a placeholder for objects' location embeddings. During the text generation, the `<trigger>` triggers the location decoding, where the hidden states of `<trigger>` can be used for both detection and segmentation, as depicted in Fig. 1. Then, the predicted or provided object location will be encoded into the embedding of the `<loc>` token for object referring.

In addition to supporting flexible output formats, the pixel2emb modeling also allows for the use of existing localization practices. While the pixel2seq paradigm can only frame the detection task as a token classification problem, the embedding-based paradigm formulates the localization task as a regression problem, enabling the adoption of established practices such as L1 loss, IoU loss and GIoU loss. In our exploration experiments, we show that our pixel2emb method enables a faster convergence and better performance for both location input (72.4 vs. 71.2) and output (79.4 vs. 71.4) tasks under the fair comparison.

Building upon the proposed pixel2emb method, we introduce a new LMM named NExT-Chat. NExT-Chat is designed to handle various conversation scenarios, including visual grounding (Fig. 4), region caption (Fig. 6), grounded image caption (Fig. 7). Thanks to the incorporation of LLM, the NExT-Chat is also capable of handling scenarios that requires grounded reasoning. By providing an extensive array of examples, we effectively demonstrate NExT-Chat's remarkable proficiency in understanding various components, including background elements, minute objects, and associating the objects with related knowledge.

To summarize, our contributions can be listed as follows:

- *Effective Method.* We propose the pixel2emb method, which can accommodate different output formats such as bounding boxes and segmentation masks.
- *NExT-Chat Model.* Based on the proposed pixel2seq method, we build NExT-Chat. To the best of our

knowledge, this is the first try to unify the chat, region input, detection and segmentation in a single LMM.

- *Experiments and Demos.* We conduct comprehensive experiments and provide abundant qualitative results to showcase the effectiveness of our proposed method.

2. Related Works

2.1. LMM

Large multimodal models (LMMs) are typically built on large language models (LLMs) and equipped with visual perception modules to enable the multimodal perception ability, which can generate captions or answer questions based on the given multimodal content. Flamingo [1] tries to extract vision information by a pre-trained vision backbone with a resampler, and incorporate them into the text features with a cross-attention mechanism. Instead of using cross-attention layers, BLIP-2 [12] and Kosmos [6] directly feed the visual features into the LLMs as soft prompts. Following BLIP-2, MiniGPT-4 [29] and VPGTrans [27] build LMMs with transfer learning, and significantly reduce the training cost. For example, VPGTrans can use only around 10% GPU hours with non-degenerated performances compared with training a new LMM from scratch. When considering the training paradigm, researchers find that a small scale instruction tuning can better align the LMM with the expected output format. MiniGPT-4 [29] fine-tunes its model with less than 5,000 self-instruct image-text pairs and turns the model into better conversation robot. Different from MiniGPT-4's self-instruct, LLaVA [14] generate the instruction tuning data with the text-only GPT-4 models by feeding the visual information as text sentences. Otter [10, 11] further propose a MIMIC-IT dataset that can turn the LMM into better in-context learners. However, these LMMs can only take the whole image/video as input and output text, and are incapable of handling region understanding tasks.

2.2. LMM for Region Reasoning

GPT4ROI [28] proposes to encode the regions as features and thus can accept the region as input. Pix2seq [5] first propose to represent object bounding box coordinates as text tokens and thus the language model can output the object locations in a token classification manner. However, pix2seq only validate its idea on traditional object detection tasks. UniTab [23] and PEVL [24] further extend the idea to vision&language tasks like visual grounding [16, 25]. Following this line, Vision-LLM [22] and Kosmos-2 [17] recently applies the token classification concept to LMMs. Take Kosmos-2 as an example, it discretize the whole image into 32×32 bins, which can be used to represent the points lying in it. Additional 32×32 tokens are introduced to the LLM's vocabulary for either

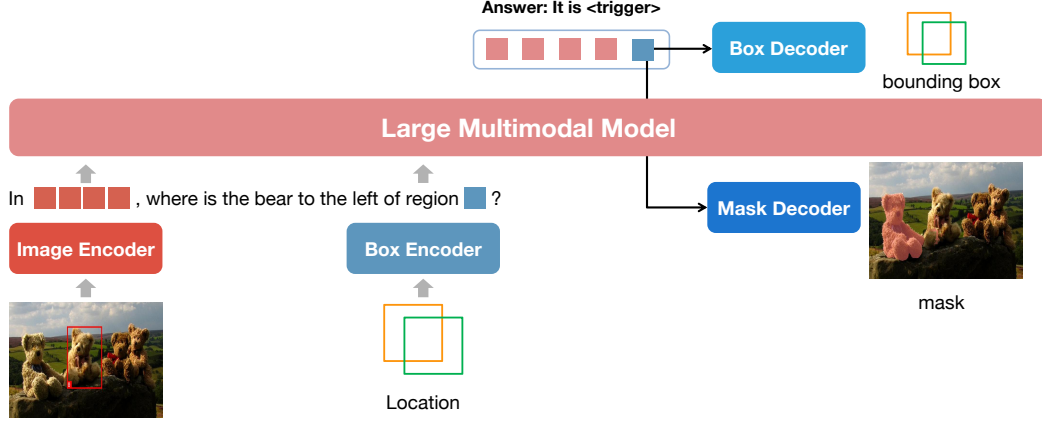


Figure 2. The overall framework of NExT-Chat. The image and given bounding boxes are encoded by image and box encoders respectively. Subsequently, the hidden states of the `<trigger>` are fed into box and mask decoders, enabling object detection and segmentation.

coordinates input or output. Thus, the LMM can achieve the region-level reasoning. Shikra [4] point out that introducing too much new tokens will inevitably increase the training difficulties. Thus, Shikra propose to reuse the LLM’s original vocabulary and turn the box coordinates into normalized numerical values with certain precision like $[0.111, 0.111, 0.333, 0.333]$. Although avoiding introducing too much new tokens, it requires roughly 26 tokens to represent each bounding box, which is ineffective. Different from these works, we do not formulate the object localization problem as a token classification problem. Our NExT-Chat introduces an `<trigger>` token as the trigger for location decoding, and then use the hidden states to decode the bounding boxes and the segmentation masks.

3. Method

In this section, we present the NExT-Chat framework, starting with an introduction to the overall LMM architecture (§3.1), followed by a description of the pixel2emb method (§3.2). Additionally, we provide details on the training process (§3.3).

3.1. LMM Architecture

For the LMM architecture, we adopt a LLaVA-like architecture. Specifically, we employ a CLIP ViT-L/14 [19] as the vision encoder. The input image is converted into 16×16 patch embeddings and then projected to the same dimension as the word embeddings of the LLM. These patch embeddings serve as visual tokens. Then, the visual tokens will be fed into a decoder-only LLM for conditional text generation. Regarding the selection of LLMs, we opt for the recently released LLaMA-2 model [21] (7B-Chat and 13B-Chat).

3.2. Pixel2Emb Method

Detection. To model the object location as output, we introduce a special token, denoted as `<trigger>`, which serves to trigger the bounding box coordinates decoding. As depicted in Fig. 2, the LMM is trained to generate the `<trigger>` token before predicting the locations. Then, the embedding $\mathbf{t} \in \mathcal{R}^n$ of `<trigger>` is then passed to the *Box Decoder* \mathcal{F} for regression. Mathematically, this can be expressed as follows:

$$\mathbf{b} = \mathcal{F}(\mathbf{t}), \quad (1)$$

where $\mathbf{b} \in \mathcal{R}^4$ represents the predicted bounding box coordinates in the format $[x_0, y_0, x_1, y_1]$.

In our NExT-Chat model, the box decoder consists of a 2-layer MLP. To supervise the location output, we employ a joint loss function comprising of the L1 loss and the GIoU loss [20] during training:

$$\mathcal{L}_{det} = \mathcal{L}_1(\mathbf{b}, \mathbf{b}_{gt}) + \alpha \text{GIoU}(\mathbf{b}, \mathbf{b}_{gt}), \quad (2)$$

where \mathbf{b}_{gt} represents the ground truth coordinates, and the ratio $\alpha = 0.4$ follows the approach utilized in DETR [3].

Segmentation. Similar to the detection process, we utilize the hidden states \mathbf{t} of the `<trigger>` as input for the mask head. Inspired by LISA [9], we use SAM [7] as our mask head, which also additionally takes the original image as input. To ensure compatibility between the hidden states and SAM, we first project the hidden states to match the dimension of SAM’s prompt embedding using a linear projector. Subsequently, the projected hidden states are fed as the prompt embedding to SAM. For improved performance, we also encode the detected bounding boxes into a prompt embedding with SAM’s prompt encoder and concatenate it with the projected embedding. To train the mask output, we

Table 1. Comparison between our P_{emb} and existing localization modeling methods. **Dec. Type** is the way of decoding coordinates. **Format** is the bounding box format. **ToB** represents the number of tokens to represent a single bounding box. Note that “[” and “]” also require two tokens. **Vocab.** is the number of new tokens added to vocabulary for location data modeling. **Rep. Models** are representative models for the given location modeling method.

Methods	Dec. Type	Format	ToB	Vocab.	Rep. Models
P_{Abin}	classification	$[x_0, y_0, x_1, y_1]$	6	224	Pix2seq, VisionLLM
P_{2bin}	classification	$[p_0, p_1]$	4	1024	Kosmos-2
P_{num}	classification	$[x_0, y_0, x_1, y_1]$	26	0	Shikra
P_{emb}	regression	<trigger> & <loc>	2	2	NExT-Chat

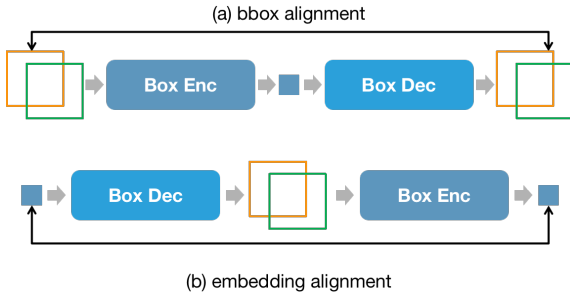


Figure 3. Cycle loss utilized to bind box encoder and decoder training.

follow the practice outlined in lightning-SAM¹:

$$\mathcal{L}_{seg} = \text{IoU}(\mathbf{m}, \mathbf{m}_{gt}) + \text{D}(\mathbf{m}, \mathbf{m}_{gt}) + \beta \text{F}(\mathbf{m}, \mathbf{m}_{gt}), \quad (3)$$

where IoU, D, and F are IoU Loss, Dice Loss, and Focal Loss separately. β is set to 20 in our experiments.

Location as Input. In addition to the location output, it is essential to incorporate location as input as well. To be consistent with the location output modeling, we also use a single embedding to represent the location information. Therefore, the output location embedding can also serve as the input embedding. Consequently, we introduce another 2-layer MLP, referred to as the location encoder \mathcal{G} . In order to simplify the problem, we convert all location formats into bounding boxes b and subsequently transform them into embeddings $\mathbf{t} \in \mathbb{R}^n$ suitable for the LLM. The location encoder can be supervised through the standard text generation loss \mathcal{L}_{text} . For instance, when inquiring about the relationship between bounding box \mathbf{b}_1 and \mathbf{b}_2 , the location encoder is compelled to provide precise information.

However, we observe that the location encoder cannot be effectively trained solely through indirect supervision from \mathcal{L}_{text} . As a result, we introduce an additional cycle loss to facilitate the training of the encoder in conjunction with the

decoder. As illustrated in Fig. 3 (a), a bounding box will be encoded and then decoded, where two bounding boxes are asked to be the same. Similarly, the hidden states of <trigger> will also be used to calculate the cycle loss (Fig. 3 (b)). Formally, the \mathcal{L}_{cyc} is defined as:

$$\mathcal{L}_{cyc} = \mathcal{L}_1(\mathbf{b}, \mathcal{F}(\mathcal{G}(\mathbf{b}))) + \mathcal{L}_2(\mathbf{t}, \mathcal{G}(\mathcal{F}(\mathbf{t}))), \quad (4)$$

where \mathbf{b} and \mathbf{t} are provided bounding box and predicted embedding respectively. Additionally, \mathcal{L}_1 and \mathcal{L}_2 correspond to the L1 Loss and L2 Loss, respectively.

3.3. Training Process

We employ a three-stage training process, consisting of pre-training, instruction tuning, and segmentation training, to train our model. The idea is to train the bounding box decoding ability for the first two stages and then extend to segmentation with a lightweight training.

Stage 1. During this stage, we perform pre-training using a mixture of data from various sources, including Flickr30K Entities [18], Visual Genome [8], RefCOCO [25], RefCOCO+ [25], RefCOCOG [16], VQA2 [2], PointQA [15], Visual7W [30], VCR [26], and LLaVA pre-training data [14]. The model is trained with a batch size of 128 and a learning rate of 2e-5 for 50k steps. During this pre-training stage, the entire language model with the box decoder, is trained while keeping the image encoder frozen. The training loss is formulated as:

$$\mathcal{L}_{s1} = \mathcal{L}_{text} + \mathcal{L}_{det} + \mathcal{L}_{cyc}. \quad (5)$$

For NExT-Chat 13B model, the stage-1 training uses 8 A100 (80G) GPUs for around 90 hours.

Stage 2. In the second stage, we further fine-tune the model using data from VQA2, RefCOCO, Flickr30K Entities, LLaVA instruction data, and Shikra-RD [4]. The batch size is reduced to 64, and the learning rate is set to 2e-6. The loss is the same with stage-1’s loss. For NExT-Chat 13B model, the stage-2 training uses 8 A100 (80G) GPUs for less than 20 hours.

¹<https://github.com/luca-medeiros/lightning-sam/tree/main>

Table 2. Comparison between our P_{emb} and other baselines on the location output tasks.

Methods	RefCOCO			RefCOCO+			RefCOCOg	
	val	testA	testB	val	testA	testB	val	test
P_{4bin}	70.9	78.1	62.6	63.0	71.4	50.7	65.2	66.8
P_{2bin}	41.8	50.5	36.1	31.3	39.2	24.9	33.7	35.3
P_{num}	71.4	78.6	62.6	61.5	70.4	51.0	62.4	63.6
P_{emb} (ours)	79.4	84.3	73.8	69.5	76.9	58.8	69.7	70.9

Table 3. Comparison between our P_{emb} and other baselines on the location input task.

Methods	Visual7W
P_{4bin}	71.2
P_{2bin}	64.9
P_{num}	51.4
P_{emb}	72.4

Stage 3. After the two stages training, the model is equipped with the ability to engage in dialogue and perform image localization. To prevent catastrophic forgetting, we keep most of the parameters frozen during the segmentation training. Specifically, we only train the linear projector between the LMM and SAM, as well as the decoder of SAM. The loss for the stage-3 is:

$$\mathcal{L}_{s3} = \mathcal{L}_{seg}. \quad (6)$$

Thanks to the small amount of training parameters, the training can be done in 4 hours with 4 A100 (80G) GPUs. This training is performed using the referring segmentation splits of RefCOCO, RefCOCO+, and RefCOCOg datasets.

4. Experiment

In this section, we begin by conducting a rigorous evaluation to validate the effectiveness of our pixel2emb approach in a fair comparison setting. Following that, we demonstrate the potential of our NExT-Chat model by presenting a wide range of qualitative results from different scenarios. Finally, we provide quantitative results to compare the performance of our NExT-Chat model with the current SOTA methods on the POPE [13] benchmark for diagnosing object hallucination.

4.1. Exploratory Study

We begin by evaluating the efficacy of our embedding-based localization modeling in comparison to existing

pixel2seq methods, specifically for both bounding box output and input tasks.

Baselines. We consider three baselines as outlined in Table 1: (1) P_{bin4} represents bounding boxes using four bins in the format of $[x_0, y_0, x_1, y_1]$. Each bin corresponds to a specific location token out of the total 224 tokens available. (2) P_{bin2} employs two bins, with each bin representing a point. The entire image is divided into 1024 discrete bins, and each bin represents the points within it. (3) P_{num} does not introduce any new tokens to the vocabulary and directly uses the textual representation of numerical values with three decimal places. In contrast, our proposed approach, denoted as P_{emb} , treats location decoding as a regression problem and introduces a special token `<trigger>` to model it.

Experiment Setup. In order to ensure a fair comparison, we maintain consistency in the architecture, batch size, learning rate, scheduler, training epochs, and datasets across all the baselines. For the location output task, we utilize the RefCOCO, RefCOCO+, and RefCOCOg datasets, where the model is required to ground the region corresponding to the provided text query. To train the model, we combine the training splits of these three datasets, while performing evaluations on the respective validation and test splits. To evaluate on the location input task, we select the Visual7W dataset [30], which necessitates the model to select the most suitable bounding box among four candidates for answering a given question. We observed that training solely on Visual7W proved to be excessively time-consuming. To mitigate this, we initialize the weights with the final checkpoints obtained from the location output task.

Results on Location Output Tasks. Table 2 provides a comprehensive comparison between our proposed P_{emb} and the baselines P_{4bin} , P_{2bin} , and P_{num} . First, we find that our P_{emb} can achieve the best performance over all of the baselines, which shows the effectiveness of our method. One reason is that our P_{emb} does not introduce too much new tokens to the vocabulary, which will significantly decrease the training difficulty under this non-pre-training sce-

Table 4. The comparison between our NExT-Chat with current SOTA models on the POPE benchmark for object hallucination diagnosis.

Datasets	Metrics	NExT-Chat	Shikra	InstructBLIP	MiniGPT-4	LLaVA	MM-GPT	mPLUG-Owl
Random	Accuracy (\uparrow)	82.34	86.90	88.57	79.67	50.37	50.10	53.97
	Precision (\uparrow)	91.85	94.40	84.09	78.24	50.19	50.05	52.07
	Recall (\uparrow)	72.13	79.27	95.13	82.20	99.13	100.00	99.60
	F1-Score (\uparrow)	80.80	86.19	89.27	80.17	66.64	66.71	68.39
	Yes	40.50	43.26	56.57	52.53	98.77	99.90	95.63
Popular	Accuracy (\uparrow)	79.80	83.97	82.77	69.73	49.87	50.00	50.90
	Precision (\uparrow)	84.98	87.55	76.27	65.86	49.93	50.00	50.46
	Recall (\uparrow)	72.40	79.20	95.13	81.93	99.27	100.00	99.40
	F1-Score (\uparrow)	78.19	83.16	84.66	73.02	66.44	66.67	66.94
	Yes	42.60	45.23	62.37	62.20	99.40	100.00	98.57
Adversarial	Accuracy (\uparrow)	79.43	83.10	72.10	65.17	49.70	50.00	50.67
	Precision (\uparrow)	82.97	85.60	65.13	61.19	49.85	50.00	50.34
	Recall (\uparrow)	74.07	79.60	95.13	82.93	99.07	100.00	99.33
	F1-Score (\uparrow)	78.27	82.49	77.32	70.42	66.32	66.67	66.82
	Yes	44.63	46.50	73.03	67.77	99.37	100.00	98.67

nario. For example, the P_{2bin} introduce the most new tokens and achieve the worst performance (*i.e.* 10 points lower than other baselines). Furthermore, unlike the baselines, which approach location decoding as a token classification problem, our P_{emb} employs a regression loss that aligns with the inherent continuity of coordinates. By incorporating these factors, our approach enhances the accuracy and effectiveness of location decoding.

Results on Location Input Tasks. In addition to reporting the location output results, we also present the findings on the location input tasks. As depicted in Table 3, our proposed P_{emb} exhibits the best performance on Visual7W, which is a region VQA task. These results demonstrate that our P_{emb} is not only adept at modeling location output but also proficient in handling location input tasks. An intriguing observation is that P_{num} , which achieves exceptional performance in Table 2, performs poorly in the location input task. We hypothesize that this discrepancy arises due to the complexity of the context in the location input task, where models are required to choose one box from four candidates. In this task, P_{num} needs to utilize $4 \times 26 = 104$ tokens to represent all the candidates, resulting in an overly intricate context for the reasoning process. Conversely, our proposed P_{emb} treats the entire bounding box as a single embedding, thereby circumventing this issue.

4.2. Applications across Different Scenarios

In this section, we present qualitative results that showcase the capabilities of our NExT-Chat model across various sce-

narios.

Visual Grounding. As shown in Fig. 4, we can see that our NExT-Chat accurately detects and segments the queried objects, such as the bears and the sky in the background. To ensure that our model is not biased towards specific objects, we test it with different queries to find all four bears individually. Our model successfully localizes each bear based on the given queries. Additionally, our model showcases reasoning abilities through challenging grounding problems. For instance, in Fig. 5, our model accurately localizes the remote in response to the query “Where is the object to control the TV in image?” It also localizes the boat based on understanding the given object location input.

Region Caption. To evaluate the effectiveness of our NExT-Chat model for region input, we conducted experiments where the model generates descriptions based on given bounding boxes. As depicted in Fig. 6, our model consistently produces accurate descriptions specifically tailored to the provided regions, without being influenced by the overall image content or salient regions. We observed this behavior consistently across different examples. Notably, in the second row of Fig. 6, our model demonstrates the ability to accurately recognize and describe small objects such as flags, as well as background objects like trees. This demonstrates the robustness and effectiveness of our model in generating region-based captions.

Object-Referenced Captioning. Another compelling application of our NExT-Chat model is its ability to describe

images by referencing specific objects present within them. Fig. 7 demonstrates that our model can accurately identify and describe the major 2 or 3 objects in an image, effectively organizing them into coherent sentences. By incorporating object references, our model demonstrates a reduced tendency to generate captions containing non-existent objects. This highlights the model’s capability to generate more accurate and contextually grounded image descriptions.

Multi-turn&Reasoning. In addition to its demonstrated ability in single-turn and concise response generation, our NExT-Chat model also possesses the capability for multi-turn conversations and generating detailed explanations in response to given questions. As illustrated in the third example of Fig. 8, our model exhibits the ability to infer the occupation of the man in the image by analyzing contextual cues such as his uniform and the horse he is riding. This inference is supported by the model’s ability to localize relevant regions within the image. Furthermore, for each hypothesis regarding the man’s occupation, our model provides detailed descriptions of the potential duties associated with that occupation. This showcases the model’s capacity for nuanced reasoning and comprehensive explanation generation.

5. Comparison with SOTAs

In this study, we evaluate our NExT-Chat model by comparing it with current state-of-the-art (SOTA) models on the hallucination diagnose datasets. Specifically, we utilize the POPE dataset to assess the models’ capability to generate correct text responses.

5.1. Hallucination

The results, presented in Table 4, demonstrate that our NExT-Chat model exhibits competitive performance compared with existing SOTA models. Notably, our model achieves the second highest accuracy on both the Popular and Adversarial splits of the POPE dataset. These findings indicate that our NExT-Chat model is competent in generating accurate responses, thus positioning it among the top-performing models in the field.

6. Conclusion

In this paper, we present a novel location modeling method called pixel2emb, which utilizes embeddings to achieve multiple location output formats, such as bounding boxes and segmentation masks. Through comprehensive exploratory experiments, we demonstrate the effectiveness of the proposed pixel2emb method. Additionally, we train a LMM named NExT-Chat, which significantly broadens the range of application scenarios for LMMs. Our NExT-Chat exhibits the ability to handle diverse tasks, including

visual grounding, region captioning, grounded captioning, multi-turn conversation, and complex question reasoning. In the future, we will continue to enhance the model’s ability on conducting better detection and segmentation. Another promising direction is to extend the NExT-Chat model to multimodal agent which can handle complex tasks that requires region understanding.

7. Limitation

In the training procedure, our dataset primarily comprises individual image inputs, resulting in a limitation of our NExT-Chat model when it comes to handling multiple image inputs. Furthermore, the absence of sufficient training data from diverse domains hinders the model’s ability to generate accurate predictions in tasks involving medical and satellite image analysis. Lastly, it is worth noting that our model is occasionally prone to errors due to the limit of image resolution (224×224).

Author Contributions

Ao Zhang initialize the project, conduct experiments and write the main part of the paper. Liming Zhao, Chen-Wei Xie, and Yun Zheng provide valuable suggestions on the method design, model training, and paper revision. They also help to coordinate the resources for the whole project. Wei Ji help to conduct the proof reading. Tat-Seng Chua provides valuable suggestions on the paper structure and paper revision.

Visual Grounding



Figure 4

Complex Grounding

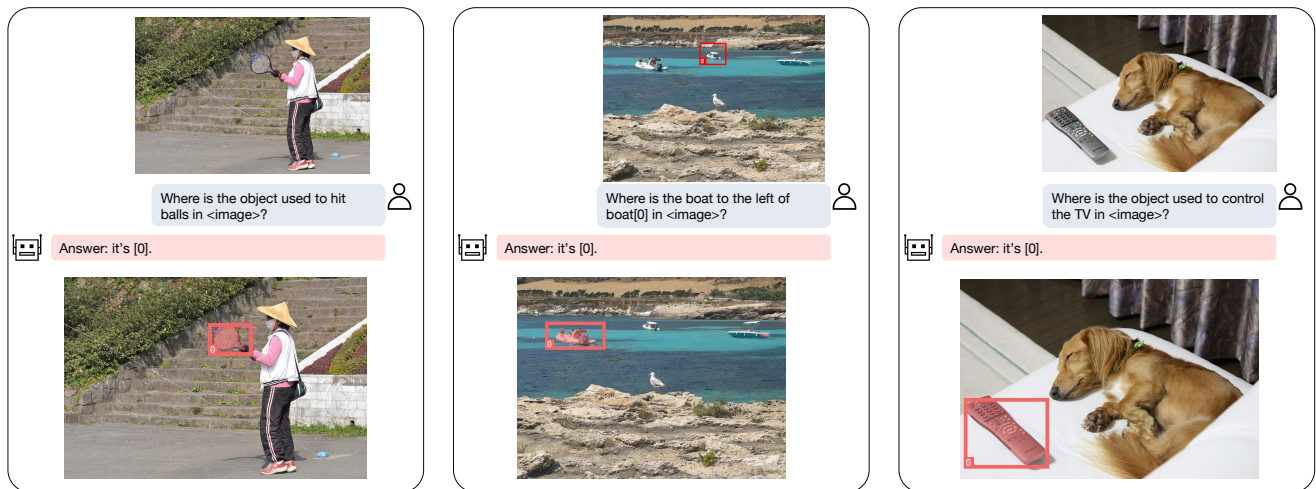


Figure 5

Region Caption

 <p>What is region [0]?</p> <p>man wearing a blue shirt.</p>	 <p>Can you describe area [0]?</p> <p>man wearing a gray shirt.</p>	 <p>Can you describe area [0]?</p> <p>stack of oranges for sale</p>
 <p>What is region [0]?</p> <p>red and white sign</p>	 <p>Can you give me a description of the region [0] in image?</p> <p>the trees are green</p>	 <p>Can you give me a description of the region [0] in image?</p> <p>a white light switch</p>

Figure 6

Grounded Caption


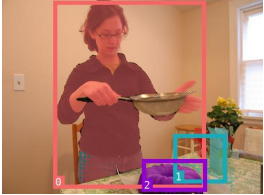

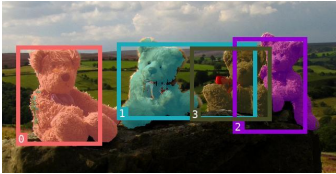



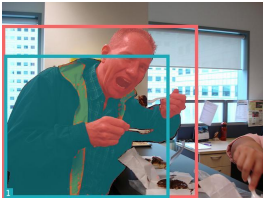

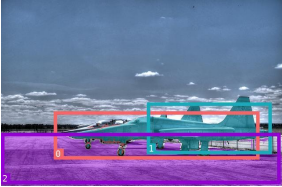



 <p>What is the woman doing? Please include object locations.</p> <p>The woman [0] is putting white powder [1] on a cake [2] .</p> 	 <p>Can you provide a description of the image <image> and include the locations for each mentioned object?</p> <p>Four stuffed teddy bears of various colors [0] [1] [2] [3] are sitting on a rock outside .</p> 	 <p>What can you tell me about the contents of the image <image>? Please indicate the positions of any mentioned objects.</p> <p>A white dog [0] is laying on a marble table top next to a cup of coffee [1] .</p> 
 <p>Describe the scene in this picture <image> and give the position info for each item you talk about.</p> <p>A man [0] in a black jacket [1] is eating .</p> 	 <p>Can you provide a description of the image <image> and include the locations for each mentioned object?</p> <p>Two fighter jets [0] [1] are parked on a runway [2] .</p> 	 <p>What can you tell me about the contents of the image <image>? Please indicate the positions of any mentioned objects.</p> <p>A brown and white dog [0] is sleeping on a white couch [1] with a television remote [2] on it .</p> 

Figure 7


Multi-turn & Reasoning




Why is the bike placed on the top of the car?

The answer is transportation.


Find the bike and share its coordinates with me.






What will the boy do in reaction to the frisbee? Please include object locations and explain.

The boy [0] will be excited and enthusiastic about the frisbee as he is a young boy who is about to catch it. He will try to catch it mid-air and will be happy about his successful catch. Once he catches the frisbee, he will throw it back to the person from whom he threw it in the first place. This will create a fun and engaging recurring play between the boy and the person. The boy will enjoy the game and the healthy competition that comes with playing frisbee.





What is the potential job of the man? Please include object locations and explain.

Based on the image, the man is wearing a uniform [0] and is riding a brown horse [1]. It is likely that the man is a police officer or serves a law enforcement role, as he is wearing a uniform, riding a horse, and possibly patrolling the area. His potential job could include patrolling neighborhoods, responding to calls or incidents, attending community events, or ensuring public safety. The man on the horse may also be a mounted patrol officer, in charge of protecting a specific area or event, watching over the public, or providing support to other law enforcement officers.

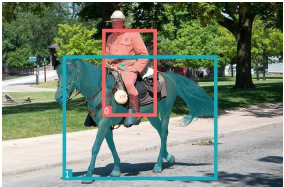


Figure 8

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. [2](#)
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015. [4](#)
- [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*, pages 213–229. Springer, 2020. [3](#)
- [4] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multi-modal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. [2](#), [3](#), [4](#)
- [5] 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. [2](#)
- [6] 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. [2](#)
- [7] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. [3](#)
- [8] 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. [4](#)
- [9] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. *arXiv preprint arXiv:2308.00692*, 2023. [3](#)
- [10] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it: Multi-modal in-context instruction tuning. 2023. [2](#)
- [11] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023. [2](#)
- [12] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. [2](#)
- [13] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023. [5](#)
- [14] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. [2](#), [4](#)
- [15] Arjun Mani, Nobline Yoo, Will Hinthorn, and Olga Russakovsky. Point and ask: Incorporating pointing into visual question answering. 2020. [4](#)
- [16] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of CVPR*, pages 11–20, 2016. [2](#), [4](#)
- [17] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. [2](#)
- [18] 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*, pages 2641–2649, 2015. [4](#)
- [19] 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*, pages 8748–8763. PMLR, 2021. [3](#)
- [20] Hamid Rezaatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 658–666, 2019. [3](#)
- [21] 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*, 2023. [3](#)
- [22] Wenhui Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, and Jifeng Dai. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023. [2](#)
- [23] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Faisal Ahmed, Zicheng Liu, Yumao Lu, and Lijuan Wang. Unitab: Unifying text and box outputs for grounded vision-language modeling. In *European Conference on Computer Vision*, pages 521–539. Springer, 2022. [2](#)
- [24] 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. [2](#)
- [25] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Proceedings of ECCV*, pages 69–85. Springer, 2016. [2](#), [4](#)

- [26] 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)*, 2019. 4
- [27] 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*, 2023. 2
- [28] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023. 2
- [29] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 2
- [30] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images, 2016. 4, 5