# Contextual Object Detection
# with Multimodal Large Language Models

**Yuhang Zang, Wei Li, Jun Han, Kaiyang Zhou, Chen Change Loy**[✉]
S-Lab, Nanyang Technological University
{zang0012, wei.l, hanj0030, kaiyang.zhou, ccloy}@ntu.edu.sg

## Abstract

Recent Multimodal Large Language Models (MLLMs) are remarkable in vision-language tasks, such as image captioning and question answering, but lack the essential perception ability, *i.e.*, object detection. In this work, we address this limitation by introducing a novel research problem of *contextual object detection*—understanding visible objects within different human-AI interactive contexts. Three representative scenarios are investigated, including the language cloze test, visual captioning, and question answering. Moreover, we present ContextDET, a unified multimodal model that is capable of end-to-end differentiable modeling of visual-language contexts, so as to locate, identify, and associate visual objects with language inputs for human-AI interaction. Our ContextDET involves three key submodels: (i) a visual encoder for extracting visual representations, (ii) a pre-trained LLM for multimodal context decoding, and (iii) a visual decoder for predicting bounding boxes given contextual object words. The new *generate-then-detect* framework enables us to detect object words within human vocabulary. Extensive experiments show the advantages of ContextDET on our proposed CODE benchmark, open-vocabulary detection, and referring image segmentation. [1]

## 1 Introduction

"*For me context is the key - from that comes the understanding of everything.*"

— Kenneth Noland

One indispensable cornerstone of computer vision—object detection—is understanding visible objects within scenes, which empowers many applications, such as robotics, autonomous driving, and AR/VR systems. Recently, Multi-modal Language Models (MLLMs) trained with internet-scale visual-language data, including Flamingo [1], PaLM-E [16], and the superb OpenAI's GPT-4 [51], have shown a revolutionary ability to allow humans to interact with AI models for various vision-language tasks, *e.g.*, image captioning and question answering. Such an interactive human-AI circumstance requires modeling *contextual* information, *i.e.*, relationships among visual objects, human words, phrases, and even dialogues. Therefore, it is desirable to advance MLLMs with the capability of locating, identifying, and associating visual objects with language inputs for human-AI interaction.

In this paper, we study a new research problem—contextual object detection—that is understanding visible objects within human-AI interactive contexts. In comparison with existing standard object detection, we consider four comprehensive objectives for such a new setting: (i) **capacity**: being able to handle a human language vocabulary; (ii) **description**: describing visual inputs from users with informative natural language statements; (iii) **perception**: locating and associating visual objects with language queries; (iv) **understanding**: complementing proper words with language hints. To

---

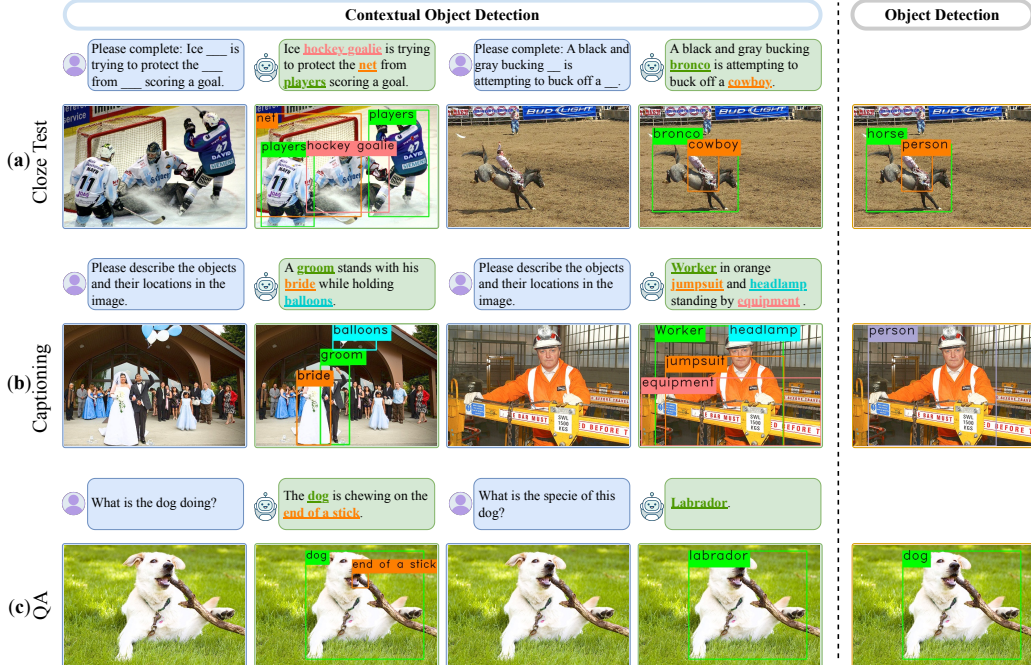[1]Github: https://github.com/yuhangzang/ContextDET.

Figure 1: We present a new **contextual object detection** task include **(a)** look at the image and complete the masked object names and locations; **(b)** predict the caption and the boxes of objects existing in the caption; **(c)** answer a question about the names and locations of objects. Unlike the traditional object detection task that typically focuses on detecting a limited set of pre-defined object classes such as 'person', our task requires predicting more specific names (*e.g.*, 'hockey goalie', 'groom', or 'bride') based on contextual understanding.

cover these four objectives, we incorporate three representative tasks: language cloze test, visual captioning, and question answering, with object detection for MLLMs (see Figure 1).

While significant progress has been made in developing more accurate and faster object detection algorithms, it is impossible to directly integrate existing deep object detectors with MLLMs for contextual object detection, due to the following reasons. First, standard deep detectors, such as Mask-RCNN [23] and DETR [6], are trained with close-set classifiers and cannot generalize well in real-world scenarios, where object categories or classes are not pre-defined or limited to a closed set. Despite the very recent development of open-vocabulary object detection [21, 87, 81, 58] that builds on state-of-the-art vision-language models (e.g., CLIP [54] and ALIGN [28]) can improve the zero-shot transfer ability for novel classes, they are constrained by the scale of pre-defined novel categories, making them incapable of detecting objects for a human language vocabulary. For example, these open-vocabulary detectors fail to handle out-of-distributed categories in Figure 1, such as hockey goalie, groom, and cowboy. Second, the inherent *locate-then-classify* paradigm of existing deep detection models is unsuitable for contextual object detection. In generic human-AI interactive scenarios, both natural objects in visual scenes and human words in language inputs have various meanings in different contexts. In Figure 1 (a) and (b), the universal 'person' category will manifest as 'goalie', 'player', 'cowboy', 'groom', 'bride', and 'worker' within distinct visual contexts. Also, as language contexts shift, the word 'labrador' supplants the representation of 'dog' (Figure 1 (c)). Consequently, an innovative detection approach is required to cater to considerably varied and changing contextual object detection.

To address the above challenges, in this work, we present ContextDET, a novel *generate-then-detect* framework, specialized for contextual object detection. Specifically, it is an end-to-end model that consists of three key modules. First, a visual encoder extracts high-level image representations for given images and produces both local and full visual tokens for further contextual modeling. Second, to effectively model multimodal contexts, we employ a pre-trained LLM to perform text generation, with conditioned inputs of both local visual tokens and task-related language tokens as the multimodal prefix. Third, taking the LLM tokens as prior knowledge for visual detection, we introduce a visual decoder that consists of multiple cross-attention layers, within which we compute conditional object queries from contextual LLM tokens, and keys and values from full visual tokens, to predict the

corresponding matching scores and bounding boxes. This allows us to detect contextual object words for a human vocabulary.

**Contributions.** In summary, our contributions are the following: (**i**) We for the first time investigate contextual object detection—a new direction for visual object detection that improves MLLMs with a greater ability for human-AI interaction. (**ii**) To open this area to empirical study, we present a new benchmark CODE with 10,346 unique object words to facilitate research on contextual object detection. (**iii**) We propose a novel *generate-then-detect* framework, ContextDET, dedicated to contextual object detection. (**iv**) We demonstrate the advantages of our ContextDET not only on the CODE benchmark but also on open-vocabulary detection and referring image segmentation tasks. We hope our work can motivate future research in contextual object detection that benefits human-AI interaction.

## 2   Related Work

**Multimodal Large Language Models (MLLMs).** Large Language Models (LLMs) have been developed to comprehend and generate textual language, showcasing remarkable performance across a wide range of Natural Language Processing (NLP) tasks. Notable examples of LLMs include OpenAI's GPT series [55, 56, 5, 50, 51], Google's T5 [57] and PaLM [11], as well as Meta's OPT [85] and LLaMA [66]. More recently, there have been advancements in the field of MLLMs [46, 67, 7, 32, 36, 26, 16, 51], exemplified by the GPT-4 model [51], which have expanded the capabilities of LLMs to comprehend both language and visual inputs. MLLMs have demonstrated impressive proficiency in a range of vision-language tasks, including image captioning and visual question answering. However, existing MLLMs are limited to generating textual outputs. In contrast, our ContextDET, built upon MLLMs, extends support to contextual object detection, providing bounding box outputs. Further comparisons are discussed in Section 4.4.

**Prompting LLMs with Vision Experts**. Several recent papers [63, 72, 77] have proposed systems that leverage the textual output generated by LLMs, such as ChatGPT [50], as prompts to manipulate external vision expert models for various vision-related tasks. In the context of object detection, these vision expert models include DETR [6], Grounding DINO [41], SAM [31], and other algorithms integrated into the HuggingFace community [27]. However, due to the frozen parameters of both LLMs and expert models, the knowledge and representations from LLMs cannot be shared, potentially leading to sub-optimal performance. In contrast to these prompting-based methods, our ContextDET employs an end-to-end training pipeline. We utilize the latent features extracted from MLLMs as conditional inputs for a visual decoder, enabling the prediction of bounding boxes.

**Object Detection with Contextual Understanding**. The term "context" commonly refers to the neighboring pixels or surrounding regions within images and has been extensively explored in previous studies to enhance object detection algorithms [14, 47, 64, 10]. In this paper, the concept of contextual information encompasses multimodal patterns and relationships, involving both visual images and textual words. Our ContextDET leverages the robust contextual understanding capability of MLLMs and applies it to the downstream object detection task. Additionally, we propose the adoption of new evaluation tasks, such as the cloze test, to more effectively assess the contextual understanding ability.

**Object Detection on Novel Classes**. Despite significant advancements in deep learning techniques [59, 42, 35, 65, 6, 9, 43, 83, 89, 70], object detection remains a challenging task in real-world scenarios, particularly in the case of zero-shot object detection [4]. Zero-shot object detection requires models trained on *base* classes to detect *novel* classes that were not encountered during training. A recent variant of zero-shot detection, known as Open-Vocabulary Object Detection, allows for the utilization of additional image-text pairs [82], garnering significant attention from the research community. In this context, recent vision and language pre-trained models [54, 86, 37, 84], such as CLIP, have been widely employed for open-vocabulary object detection [21, 87, 17, 81, 58, 33, 73, 74, 75]. Instead of relying solely on CLIP, our ContextDET demonstrates that MLLMs can also be applied effectively to the open-vocabulary setting. With the assistance of MLLMs, ContextDET is not constrained by pre-defined *base* or *novel* classes. Notably, the object names predicted by ContextDET can be generated as the most contextually valid English words by the MLLMs.

**Visual Grounding**. Visual grounding tasks, such as referring expression comprehension [30], involve combining object detection with language understanding abilities. In these tasks, a language query
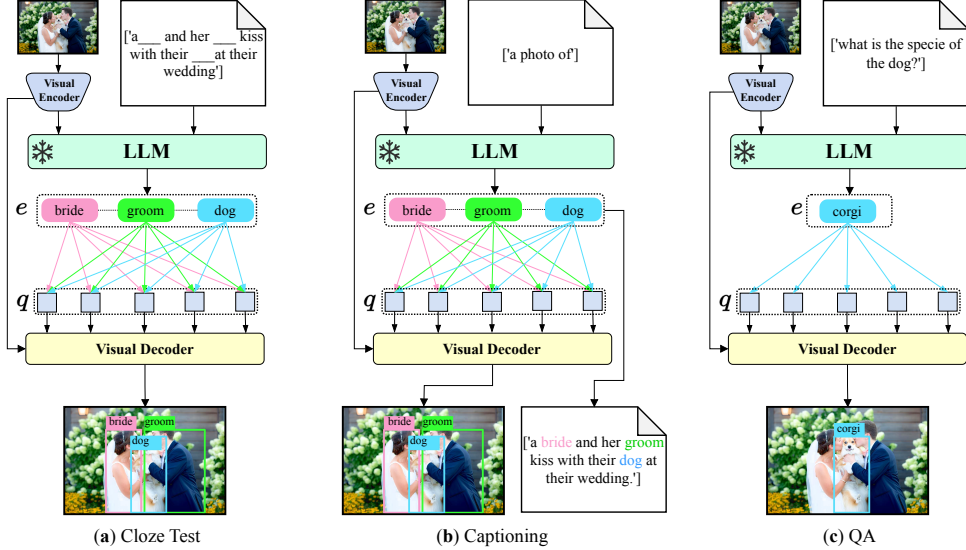
**Figure 2:** Our ContextDET is a unified end-to-end framework, being capable of taking different language token inputs for different tasks, including **(a)** cloze test **(b)** captioning and **(c)** question answering. ❄: frozen. The symbol $e$ indicates latent embeddings of LLM (Section 3.2), and the symbol $q$ denotes object queries of the visual decoder (Section 3.3).

is provided to describe a specific object, and models are tasked with predicting the position of the referred object. State-of-the-art algorithms [76, 71] commonly employ Transformer-based cross-modal structures or multimodal pre-training [29]. Our proposed contextual object detection task presents even greater challenges compared to visual grounding. For example, in our cloze test, the language query is incomplete, and the object names are masked. Models are required to infer both the missing object name words and their positions based on contextual information. Furthermore, in our contextual captioning setting, no language query is given. Additionally, in our contextual QA setting, the objects are described using human language in an *interactive* environment.

## 3 Approach

This section describes our contextual object detection framework, ContextDET, which accepts images interleaved with human text as inputs and produces free-form text and corresponding bounding boxes as outputs. As illustrated in Figure 2, our ContextDET is end-to-end and consists of three key architectural components: (1) a visual encoder that extracts high-level image representations and computes visual tokens, (2) a pre-trained LLM that decodes multimodal contextual tokens with a task-related multimodal prefix, and (3) a visual decoder that predicts matching scores and bounding boxes for conditional queries linked to contextual object words.

### 3.1 Visual Encoder

Given an image input $x \in \mathbb{R}^{3 \times H \times W}$, we use a vision backbone parameterized by $\phi$ to extract image-level spatial features $v = f_\phi(x) \in \mathbb{R}^{d \times h \times w}$, where $d$ denotes the feature dimension. The vision backbone $\phi$ is pre-trained and frozen, which can be selected from various options, including ResNet [24], Vision Transformer (ViT) [15], or Swin Transformer [43]. Subsequently, the image-level features $v$ are transformed into two distinct representations.

**Local Visual Tokens.** We first divide the 2D spatial grid of features as $p$ local bins and apply adaptive average pooling for each bin, followed by a linear projection then flattened to 1D: $z = \texttt{Linear}(\texttt{AvgPool}(v))$. As a result, fixed-sized visual tokens $z \in \mathbb{R}^{d_1 \times p}$ are obtained and fed to the LLM (Section 3.2), Here, $d_1$ represents the input dimension of LLM.

**Full Visual Tokens.** We flatten the 2D spatial features $v$ as 1D sequence with $m = h \times w$ tokens and leverage six Transformer layers $\psi$ to compute the encoded full visual tokens: $c = f_\psi(v) \in \mathbb{R}^{d_2 \times m}$, which will serve as inputs for the visual decoder (Section 3.3).

4

## 3.2 Multimodal Context Modeling with LLM

Motivated by the finding that LLMs are strong context generators [80] for solving various knowledge-intensive tasks, it is thus appealing to model multimodal contexts with LLMs. We consider performing text generation with the LLM, conditioned on both the visual representations produced by the visual encoder described in Section 3.1 and task-oriented human languages.

**Multimodal Tokens.** Given the visual context of input images, we generate language contexts that describe the visual information or complement missing words. Specifically, the inputs to the LLM consist of (1) the local visual tokens $z \in \mathbb{R}^{d_1 \times p}$, and (2) a series of language tokens $t_{1:l} = \{t_1, \ldots, t_l\} \in \mathbb{R}^{d_1 \times l}$, where the symbol $l$ is the sequence length of the language tokens. The language tokens $t_{1:l}$ have different forms for different contextual object detection settings. For the cloze test, the language tokens are tokenized sentences with masked names, *e.g.*, 'a [MASK] and her [MASK] kiss with their [MASK] at their wedding'. For the visual captioning, the language tokens are tokenized text prompts—'a photo of'—to describe the image. For the question answering, the language tokens represent the tokenized sentences of questions, *e.g.*, 'Question: what is the specie of the dog? Answer:'.

**Multimodal Prefixed LLM Decoding.** A pre-trained LLM $\theta$ can be conditioned on a prefix $w_{1:n}$ that contains multimodal tokens to generate text in an autoregressive way:

$$p(w_{n+1:L}|w_{1:n}) = \prod_{i=n+1}^{L} p_\theta(w_{i+1}|w_{1:i}). \tag{1}$$

Here, the prefix $w_{1:n} = [z, t_{1:l}] \in \mathbb{R}^{d_1 \times (p+l)}$ is obtained via concatenating the local visual tokens $z$ with a sequence of language tokens $t_{1:l}$. Specifically, the LLM consists of multiple Transformer layers (`TransLayers`) and a final Feed Forward Network (`FFN`). To generate new tokens, the LLM first predicts the latent embedding $e_{n+1}$ for the new $n+1$-th token:

$$e_{n+1} = \texttt{TransLayers}(w_{1:n}), \tag{2}$$

which contains decoded multimodal contextual information. Then, the `FFN` computes the probability distribution $p(w_{n+1})$ based on the latent embedding $e_{n+1}$:

$$p(w_{n+1}) = \texttt{Softmax}(\texttt{FFN}(e_{n+1})), \tag{3}$$

where the tokens $w_{n+1}$ are elements of a vocabulary $\mathcal{W}$ that corresponding to human words in natural language. Such autoregressive generation ends when the generated language token $w_L$ hits the [EOS] token, *i.e.*, the ending of sentences.

## 3.3 Visual Decoder

In order to associate object words with corresponding visual objects in given images, we propose a novel *generate-then-detect* pipeline for contextual object detection. Unlike the common *detect-then-classify* pipeline in standard object detectors (*e.g.*, Mask R-CNN [23] and DETR [6]) that exhaustively locate and recognize all possible objects as pre-defined categories, we consider using the LLM tokens as prior knowledge for visual detection. This allows us to detect contextual object words, while not being limited to a close set of object classes.

**Contextual LLM Tokens as Conditional Object Queries.** From both language prefix $t_{1:l}$ and generated tokens $w_{n+1:L}$ (Section 3.2), we predict the binary-classification probability of noun object words. Then, we automatically select those language tokens related to object words (*e.g.*, 'bride', 'groom', 'dog') as contextual object tokens and take their latent embeddings as conditional inputs for the visual decoder. To be specific, we set up $N$ learnable object queries $q$ as learnable positional embeddings in the visual decoder. For a contextual token, *e.g.*, 'bride', we obtain the conditional object queries that linked to 'bride', by incorporating the corresponding latent embedding $e$ from the LLM with the object queries:

$$\bar{q} = q + \texttt{Linear}(\texttt{Repeat}(e)). \tag{4}$$

Here, we repeat the latent embedding $e$ for 'bride' $N$ times so as to align with the number of the object queries $q$. Also, a linear layer is employed for dimension projection.

**Conditional Multimodal Context Decoding.** To model cross-modal contextual relationships, we employ six Transformer cross-attention layers in the visual decoder, in which the keys and values are

obtained from the full visual tokens $c$ extracted by the visual encoder (Section 3.1) while the queries are derived from the conditional object queries $\bar{q}$ for computing cross-attention:

$$\hat{q} = \texttt{CrossAttenLayers}(c, \bar{q}). \tag{5}$$

By doing so, the visual decoder learns to focus on specific areas of the visual context that are relevant to the conditional query for 'bride'.

**Box and Matching Predictions for Contextual Words.** Finally, we compute the binary matching score and box prediction from the output latent embedding $\hat{q}$ using two $\texttt{FFN}$ prediction heads:

$$p = \texttt{FFN}_{\text{cls}}(\hat{q}) \in \mathbb{R}^{N \times 2}, b = \texttt{FFN}_{\text{box}}(\hat{q}) \in \mathbb{R}^{N \times 4}, \tag{6}$$

where $p$ refers to the probability of being matched or not matched given the conditional object word, and $b$ indicates the predicted box coordinates.

**Conditional Matching for Label Assignment.** We introduce a conditional modification to the default optimal bipartite matching in DETR [29] that finds the best match between the set of $N$ predictions and the set of ground truth objects. In our approach, only the ground-truth bounding boxes that match the conditional object words are involved in the loss computation. This conditional matching ensures that the model focuses solely on the objects described by the language queries.

## 3.4 Training Details

We use multi-scale deformable attention [89] and IoU-based label assignment [52] to accelerate the convergence speed. The vision encoder $\phi$ also supports the pre-trained weights from previous MLLM such as BLIP-2 [36].

**Loss Function.** In Section 3.3, we use conditional matching to derive the label assignments, which include the ground-truth matching labels $\hat{p}$ and the associated box coordinates $\hat{b}$. For our predicted language token $w$, we can straightforwardly get the annotated ground truth token $\hat{w}$, *e.g.*, tokenized answers for the cloze test. We can also obtain the annotated binary label $\bar{w}$ indicating whether a token belongs to an object word or not. Based on the label assignment results, the overall loss function $\mathcal{L}$ is defined as:

$$\mathcal{L} = \lambda_{\text{cls}}\mathcal{L}_{\text{cls}}\left(p, \hat{p}\right) + \lambda_{\text{box}}\mathcal{L}_{\text{box}}(b, \hat{b}) + \lambda_{\text{lm}}\mathcal{L}_{\text{lm}}(w, \hat{w}) + \lambda_{\text{noun}}\mathcal{L}_{\text{noun}}(w, \bar{w}) \tag{7}$$

Here, the classification loss $\mathcal{L}_{\text{cls}}$ is a binary softmax classification loss of two classes: matched *vs*. not matched. The box-related loss $\mathcal{L}_{\text{box}}$ is either L1 loss or GIoU loss [60]. The language modeling loss $\mathcal{L}_{\text{lm}}$ is softmax classification loss over the vocabulary size $\mathcal{W}$ of the LLM Tokenizer. The noun loss $\mathcal{L}_{\text{noun}}$ is a binary classification loss that determines whether a token is an object word or not. We set the loss weighting hyper-parameters $\lambda_{\text{cls}} = 1$, $\lambda_{\text{box}} = 5$, $\lambda_{\text{lm}} = 1$, and $\lambda_{\text{noun}} = 1$.

## 4 Experiments

We present the results of ContextDET on different tasks, including our proposed contextual object detection task discussed in Section 4.1, open-vocabulary object detection described in Section 4.3, and the referring image segmentation task presented in the Appendix. For contextual object detection, we focus on providing quantitative and qualitative results for the cloze test setting, as inferring relevant object words from vast human vocabulary poses a significant challenge. Additionally, we provide qualitative results for both contextual captioning and contextual question-answering settings.

**Implementation Details.** Our proposed method is implemented in PyTorch and all models are trained using a single machine with 4 NVIDIA A100 GPUs. During training, data augmentation techniques are applied including random horizontal flipping with a probability of 0.5 and large-scale jittering [19]. We set the batch size to 8 and train the model for 6 epochs. We use AdamW [44] optimizer with a learning rate of $1e^{-4}$ and a weight decay of $0.05$.

## 4.1 Contextual Object Detection

**CODE Dataset.** To facilitate research on contextual object detection, we construct a Contextual Object DEtection (CODE) dataset. Specifically, we collected images, bounding boxes and captions

Table 1: Benchmark results of ContextDET on our CODE dataset `val` set. We report four metrics for comparisons: Acc@1, Acc@5, AP@1, and AP@5. We also report the total number of parameters, number of trainable parameters, training time $T_{train}$ and testing time $T_{test}$ for efficiency analysis.

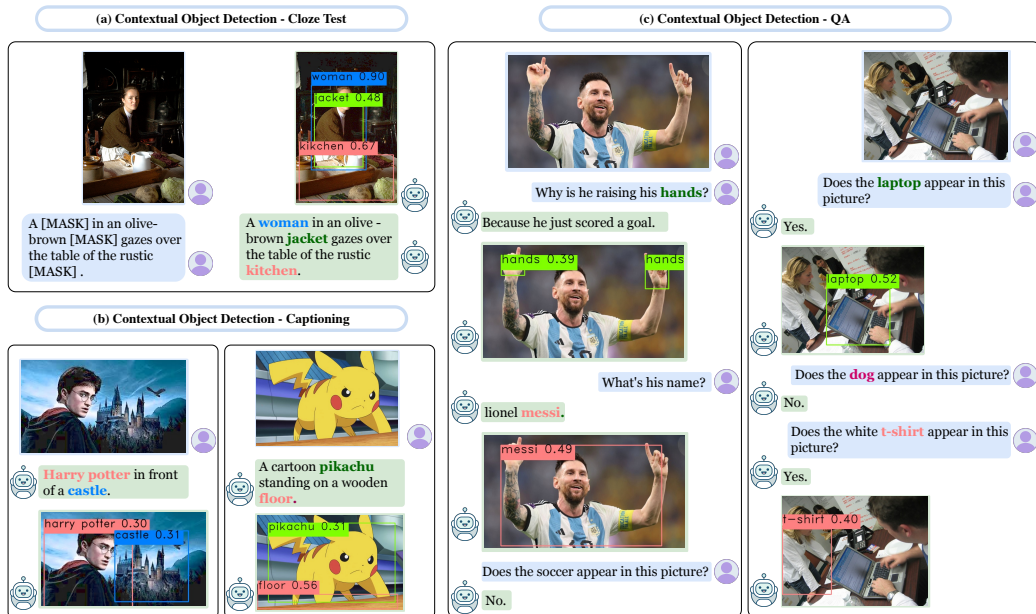| # | Language Model | Vision Backbone | Acc@1 | Acc@5 | AP@1 | AP@5 | Total #Params (M) | Learnable #Params (M) (%) | $T_{train}$ (s/iter) | $T_{test}$ (s/iter) |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | OPT-2.7B | ResNet50 | 48.7 | 73.8 | 10.2 | 20.5 | 2835 | 183 | 0.437 | 0.224 |
| 2 | | Swin-B | 54.3 | 78.1 | 13.1 | 25.3 | 2893 | 241 | 0.625 | 0.241 |
| 3 | OPT-6.7B | ResNet50 | 49.2 | 74.5 | 11.1 | 23.4 | 6922 | 263 | 0.448 | 0.248 |
| 4 | | Swin-B | **54.8** | **78.6** | **13.7** | **26.6** | 6979 | 320 | 0.652 | 0.251 |



Figure 3: Qualitative examples predicted by ContextDET in our three contextual object detection settings include (a) cloze test, (b) captioning, and (c) question answering. The 'harry potter', 'pikachu', and 'messi' are novel names that are not annotated in the CODE training set. ContextDET shows plausible contextual understanding and generalization abilities.

annotations from Flickr30k [78] and Flickr30k Entities [53]. We added annotations containing the position information of object names in the caption strings. These object names will be replaced with '[MASK]' tokens to serve as input in our *cloze test* setting. CODE is divided into three splits: the `train` split has 665,161 bounding boxes in 29,781 images, the `val` split has 22,061 bounding boxes in 1,000 images, and the `test` split has 21,641 bounding boxes in 999 images. In total, the CODE dataset has 10,346 unique object names, surpassing the number of object names in any previous detection dataset, such as COCO [38] (80 classes) and LVIS [22] (1,203 classes). Please refer to Appendix for more details about our CODE dataset.

**Evaluation Metrics.** In our contextual cloze test setting, we compute accuracy by calculating the percentage of correctly predicted object words. However, evaluating this accuracy poses a challenge due to the presence of numerous synonyms and fine-grained object words in human language, which can be difficult for annotators to distinguish. This is a problem similar to those faced by previous large vocabulary image-classification datasets, such as ImageNet [12], which use the top-5 accuracy metric as a supplementary metric to the top-1 accuracy. Consequently, we also adopt both the top-1 accuracy (Acc@1) and the top-5 accuracy (Acc@5) as our evaluation metrics. For box evaluation, we compute the mean Average Precision (mAP) metric based on the top-1 and top-5 predicted names, which are represented as AP@1 and AP@5. In evaluation, we compared the object name words rather than pre-defined category IDs, which allows a flexible extension to accommodate a vast human vocabulary. The Appendix contains more implementation details.

**Results.** We provide the results of ContextDET on the CODE dataset in Table 1. We first report the results using OPT-2.7B [85] as the language model and ResNet50 [24] as the vision backbone (Row #1). Our results suggest that the contextual cloze test task is very challenging: the top-1 AP

7

Table 2: Ablation studies on the impact of using local visual tokens $z$.

| # | $z$ | Acc@1 | Acc@5 | AP@1 | AP@5 |
|---|---|---|---|---|---|
| 1 | ✗ | 30.9 | 57.1 | 4.0 | 13.6 |
| 2 | ✓ | **48.7** | **73.9** | **10.4** | **21.6** |

Table 3: Ablation study: varying values of $p$.

| # | $p$ | Acc@1 | Acc@5 | AP@1 | AP@5 |
|---|---|---|---|---|---|
| 1 | 4 | 48.4 | 73.2 | 10.1 | 20.1 |
| 2 | 9 | **48.7** | **73.9** | **10.4** | **21.6** |
| 3 | 16 | 47.5 | 72.9 | 9.9 | 19.4 |

(AP@1) is just 10.2, which falls significantly below the performance of previous object detection datasets like COCO. Moreover, our study suggests that using more powerful language models and vision backbones can improve performance. When we replace ResNet50 with Swin-B [43] (Row #2), we observe a notable improvement from 10.2 to 13.1 in AP@1. In addition, by replacing OPT-2.7B with the larger OPT-6.7B (Row #4), we achieve an even higher AP@1 performance of 13.7.

**Efficiency Analysis.** The majority of parameters in our model, including the LLM component, are frozen, resulting in a small percentage of learnable parameters. As shown in Table 1 Row #1, when employing OPT-2.7B and the ResNet50 backbone, only 6.4% (183 out of 2,835) of parameters are trainable. Our design does not impose a significant computational burden and can be easily reproduced.

**Visualization.** Besides the quantitative evaluation on the CODE benchmark, we further qualitatively evaluate ContextDET using more diverse images and objects, as shown in Figure 3. We observe the capacity of ContextDET for complex contextual understanding and generalization to open-world names. For example, as illustrated in Figure 3 (a), ContextDET can reasonably infer the object names to fill the masked tokens, and accurately connect the object names with bounding boxes. Moreover, ContextDET is capable of predicting the names and locations of open-world concepts (*e.g.*, 'Harry Potter', 'Pikachu', 'Messi'), which are difficult to detect using previous close-set object detectors. Finally, in Figure 3 (c), we show that ContextDET can engage in multi-round question-answering conversations, and predict the bounding boxes of objects mentioned in the dialog history. Please refer to the appendix for more qualitative examples including failure cases.

## 4.2 Ablation Studies

We investigate the effects of using local visual tokens $z$ and the associated hyper-parameter $p$ that determines the number of local bins. The experiments are conducted on the CODE `val` set.

**LLM without Local Visual Tokens.** In our contextual cloze test setting, LLM is capable of making predictions even without the presence of the local visual token input $z$. However, upon analyzing the results presented in Table 2, we observe a significant performance drop. For example, the top-1 accuracy drops around 20 percent from 48.7 to 30.9 (%). This observation emphasizes the crucial role of adding visual local tokens in our method for contextual understanding. We also observe that the value of language modeling loss $\mathcal{L}_{lm}$ barely decreases in the absence of $z$.

**Hyper-Parameter $p$.** As discussed in Section 3.1, we have $p$ visual local tokens that serve as prefix inputs for LLM decoding. In Table 3, we show the effects of using different values for $p$. We observe that selecting $p = 9$ (Row #2) yields the optimal results, making it our default choice.

## 4.3 Open-Vocabulary Object Detection

We demonstrate that our proposed ContextDET can also be applied to the open-vocabulary object detection task, aiming to evaluate the generalization ability. Following previous works [4, 82], we use the OV-COCO benchmark and divide 65 categories as the split of base/novel (48/17) classes. The model is trained on the base classes only but evaluated on the novel classes (unavailable during model training). We measure the performance with the Average Precision (AP) metric on the base, novel, and all classes.

To adapt ContextDET into the open-vocabulary setting, we ask questions like 'Does the [CLASS] appear in this picture?' for every class including base and novel classes. If MLLM responds with a positive answer 'Yes', we take the latent embedding $e$ of the corresponding class name as a conditional input for our visual decoder (Section 3.3). We compare ContextDET with selected baseline methods including the state-of-the-art method BARON [74] in Table 4. We observe that ContextDET significantly outperforms BARON by large margins of 2.8%, 4.7%, and 4.2% on the novel, base, and all sets, respectively. All the baseline methods rely on prior knowledge

Table 4: Comparison with state-of-the-art open-vocabulary detection methods on OV-COCO benchmark.

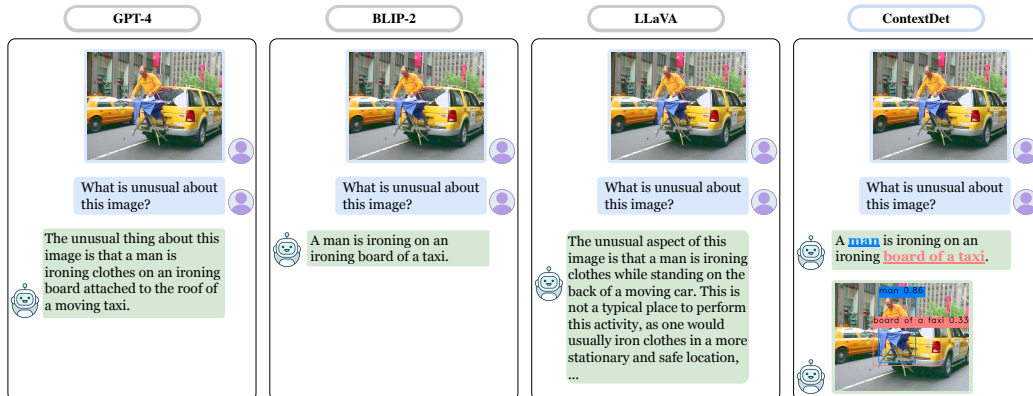| # | Method | Venue | CLIP | MLLM | Backbone | $AP_{50}^{novel}$ | $AP_{50}^{base}$ | $AP_{50}$ |
|---|--------|-------|------|------|----------|------|------|------|
| 1 | ViLD [21] | ICLR'22 | ✓ | ✗ | ResNet50-FPN | 27.6 | 59.5 | 51.2 |
| 2 | OV-DETR [81] | ECCV'22 | ✓ | ✗ | ResNet50 | 29.4 | 61.0 | 52.7 |
| 2 | BARON [74] | CVPR'23 | ✓ | ✗ | ResNet50-FPN | 34.0 | 60.4 | 53.5 |
| 4 | ContextDET | - | ✗ | ✓ | ResNet50 | **36.8** | **65.1** | **57.7** |



Figure 4: Qualitative examples comparing ContextDET with existing Multimodal Language Models (MLLMs), including GPT-4 [51], BLIP-2 [36], and LLaVA [39]. Our method predicts related bounding boxes for the object names mentioned in the text outputs, (*e.g.*, 'man', 'board of a taxi'), enabling a more comprehensive interpretation for visual-language tasks and paving the way for broader application areas.

from the vision-language model CLIP. In contrast, our ContextDET uses MLLM to detect novel objects. The results show that MLLM trained on web-scale datasets has strong generalizability that could benefit the open-vocabulary task.

### 4.4 Qualitative Results

**Comparison with MLLMs.** We present some visual examples in Figure 4 and compare our ContextDET with some popular MLLMs like GPT-4 [51]. Existing MLLMs can only generate textual outputs while our ContextDET pushes the boundaries further by providing bounding boxes of objects of interest. In particular, our method allows fine-grained localization of objects of interest specified in the text input, which offers a higher degree of interpretability for vision-language models. Broadly speaking, our method offers new possibilities for various applications requiring both object localization and conversational interaction, *e.g.*, AR/VR systems and robotics.

## 5    Conclusion

Although recent MLLMs have demonstrated remarkable abilities in vision-language tasks such as question-answering, their potential in other perception tasks remains largely unexplored. Our ContextDET highlights the significant potential of MLLMs in diverse perception tasks, such as the proposed contextual object detection task, which predicts precise object names and their locations in images for human-AI interaction. To train our model, we needed to associate object words of bounding boxes with language captions, incurring a high annotation cost. Consequently, we used less training data compared to previous MLLM papers, which may limit our final performance. In future work, we plan to explore the use of semi-supervised or weakly-supervised learning techniques to reduce annotation costs. Additionally, apart from their contextual understanding ability, we believe that other abilities of MLLMs remain underexplored for downstream tasks, such as their interactive ability for instruction tuning. For instance, can MLLMs be utilized to post-process detection outputs based on human language instructions? By providing instructions such as "shift the predicted box slightly to the left," "remove the redundant overlapped box," or "correct the predicted class from eagle to falcon," can MLLMs adjust the predictions accordingly to meet our expectations? We hope the insights presented in this paper could inspire further research into adapting MLLMs to revolutionize more computer vision tasks.

# References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022.

[2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*, 2018.

[3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. VQA: Visual question answering. In *ICCV*, 2015.

[4] Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot object detection. In *ECCV*, 2018.

[5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020.

[6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.

[7] Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed Elhoseiny. VisualGPT: Data-efficient adaptation of pretrained language models for image captioning. In *CVPR*, 2022.

[8] Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. SCA-CNN: Spatial and channel-wise attention in convolutional networks for image captioning. In *CVPR*, 2017.

[9] Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2Seq: A language modeling framework for object detection. In *ICLR*, 2022.

[10] Zhe Chen, Shaoli Huang, and Dacheng Tao. Context refinement for object detection. In *ECCV*, 2018.

[11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. PaLM: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.

[13] Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. Vision-language transformer and query generation for referring segmentation. In *ICCV*, 2021.

[14] Santosh K Divvala, Derek Hoiem, James H Hays, Alexei A Efros, and Martial Hebert. An empirical study of context in object detection. In *CVPR*, 2009.

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.

[16] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. PaLM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.

[17] Yu Du, Fangyun Wei, Zihe Zhang, Miaojing Shi, Yue Gao, and Guoqi Li. Learning to prompt for open-vocabulary object detection with vision-language model. In *CVPR*, 2022.

[18] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. *IJCV*, 2010.

[19] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *CVPR*, 2021.

[20] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017.

[21] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *ICLR*, 2022.

[22] Agrim Gupta, Piotr Dollar, and Ross Girshick. LVIS: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019.

[23] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *CVPR*, 2017.

[24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[25] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 1997.

[26] 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.

[27] HuggingFace. Huggingface. `https://huggingface.co/`.

[28] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021.

[29] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *CVPR*, 2021.

[30] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015.

[31] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.

[32] Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for multi-modal generation. *arXiv preprint arXiv:2301.13823*, 2023.

[33] Weicheng Kuo, Yin Cui, Xiuye Gu, AJ Piergiovanni, and Anelia Angelova. F-VLM: Open-vocabulary object detection upon frozen vision and language models. *arXiv preprint arXiv:2209.15639*, 2022.

[34] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, et al. The open images dataset v4. *IJCV*, 2020.

[35] Hei Law and Jia Deng. CornerNet: Detecting objects as paired keypoints. In *ECCV*, 2018.

[36] 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.

[37] Liunian Harold Li*, Pengchuan Zhang*, Haotian Zhang*, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded language-image pre-training. In *CVPR*, 2022.

[38] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.

[39] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.

[40] Jiang Liu, Hui Ding, Zhaowei Cai, Yuting Zhang, Ravi Kumar Satzoda, Vijay Mahadevan, and R Manmatha. PolyFormer: Referring image segmentation as sequential polygon generation. In *CVPR*, 2023.

[41] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding DINO: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.

[42] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. SSD: Single shot multibox detector. In *ECCV*, 2016.

[43] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical vision transformer using shifted windows. In *CVPR*, 2021.

[44] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.

[45] Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. DetectGPT: Zero-shot machine-generated text detection using probability curvature. *arXiv preprint arXiv:2301.11305*, 2023.

[46] Ron Mokady, Amir Hertz, and Amit H Bermano. ClipCap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021.

[47] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, 2014.

[48] Li Muchen and Sigal Leonid. Referring Transformer: A one-step approach to multi-task visual grounding. In *NeurIPS*, 2021.

[49] Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. Modeling context between objects for referring expression understanding. In *ECCV*, 2016.

[50] OpenAI. Chatgpt: Optimizing language models for dialogue, 2022. `https://openai.com/blog/chatgpt`.

[51] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[52] Jeffrey Ouyang-Zhang, Jang Hyun Cho, Xingyi Zhou, and Philipp Krähenbühl. NMS strikes back. *arXiv preprint arXiv:2212.06137*, 2022.

[53] 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 *CVPR*, 2015.

[54] 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 *ICML*, 2021.

[55] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. *OpenAI blog*, 2018.

[56] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.

[57] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 2020.

[58] Hanoona Rasheed, Muhammad Maaz, Muhammad Uzair Khattak, Salman Khan, and Fahad Shahbaz Khan. Bridging the gap between object and image-level representations for open-vocabulary detection. In *NeurIPS*, 2022.

[59] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015.

[60] Hamid Rezatofighi, 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 *CVPR*, 2019.

[61] Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi. Can ai-generated text be reliably detected? *arXiv preprint arXiv:2303.11156*, 2023.

[62] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *CVPR*, 2019.

[63] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. HuggingGPT: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*, 2023.

[64] Abhinav Shrivastava and Abhinav Gupta. Contextual priming and feedback for faster r-cnn. In *ECCV*, 2016.

[65] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: Fully convolutional one-stage object detection. In *CVPR*, 2019.

[66] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[67] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. In *NeurIPS*, 2021.

[68] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.

[69] Jiaqi Wang, Pan Zhang, Tao Chu, Yuhang Cao, Yujie Zhou, Tong Wu, Bin Wang, Conghui He, and Dahua Lin. V3Det: Vast vocabulary visual detection dataset. *arXiv preprint arXiv:2304.03752*, 2023.

[70] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. InternImage: Exploring large-scale vision foundation models with deformable convolutions. In *CVPR*, 2023.

[71] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. CRIS: Clip-driven referring image segmentation. In *CVPR*, 2022.

[72] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual ChatGPT: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.

[73] Jianzong Wu, Xiangtai Li, Henghui Ding, Xia Li, Guangliang Cheng, Yunhai Tong, and Chen Change Loy. Betrayed by captions: Joint caption grounding and generation for open vocabulary instance segmentation. *arXiv preprint arXiv:2301.00805*, 2023.

[74] Size Wu, Wenwei Zhang, Sheng Jin, Wentao Liu, and Chen Change Loy. Aligning bag of regions for open-vocabulary object detection. In *CVPR*, 2023.

[75] Xiaoshi Wu, Feng Zhu, Rui Zhao, and Hongsheng Li. CORA: Adapting clip for open-vocabulary detection with region prompting and anchor pre-matching. In *CVPR*, 2023.

[76] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. LAVT: Language-aware vision transformer for referring image segmentation. In *CVPR*, 2022.

[77] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. MM-REACT: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023.

[78] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *TACL*, 2014.

[79] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *ECCV*, 2016.

[80] Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. Generate rather than retrieve: Large language models are strong context generators. In *ICLR*, 2022.

[81] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-vocabulary detr with conditional matching. In *ECCV*, 2022.

[82] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *CVPR*, 2021.

[83] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. DINO: Detr with improved denoising anchor boxes for end-to-end object detection. In *ICLR*, 2023.

[84] Haotian* Zhang, Pengchuan* Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Harold Li, Xiyang Dai, Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. GLIPv2: Unifying localization and vision-language understanding. In *NeurIPS*, 2022.

[85] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.

[86] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luowei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. RegionCLIP: Region-based language-image pretraining. In *CVPR*, 2022.

[87] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022.

[88] Chaoyang Zhu, Yiyi Zhou, Yunhang Shen, Gen Luo, Xingjia Pan, Mingbao Lin, Chao Chen, Liujuan Cao, Xiaoshuai Sun, and Rongrong Ji. SeqTR: A simple yet universal network for visual grounding. In *ECCV*, 2022.

[89] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable transformers for end-to-end object detection. In *ICLR*, 2021.

# Appendix

In this supplementary, we discuss more related works, dataset and evaluation details, more experimental results, and the broader impact of this work.

- In Section A, we present more discussion about related works.
- In Section B, we discuss more details about our proposed CODE benchmark.
- In Section C, we provide the evaluation details on our proposed contextual cloze test setting.
- In Section D, we show more experiments on the referring image segmentation task and qualitative results on the contextual object detection task.
- In Section E, we provide the discussion about the broader impact of this paper.

## A    More Related Works

In this section, we discuss more related tasks that are not fully covered in the main paper, including image captioning and visual question answering. Table 5 also summarizes the differences between our proposed three contextual object detection settings with previous related tasks.

**Image Captioning.** Image captioning focuses on generating descriptive sentences to understand given images. Typically, image captioning models first encode the input image as feature embeddings using pre-trained classification [8], object detection [2] or vision language models [46]. Subsequently, submodules like LSTMs [25] or Transformers [68] are employed to decode feature embeddings into predicted sentences. In contrast, our *contextual captioning* task extends beyond language outputs by requiring the model to predict the locations of the bounding boxes containing the objects mentioned in the generated captions.

**Visual Question Answering (VQA).** Visual question answering tasks involve answering questions related to given images [3, 20]. In traditional VQA, model inputs and outputs are comprised of question-answer pairs in natural language. However, in our *contextual QA* task, questions are specifically focused on inquiring about object names and locations, while corresponding answers are expected to include the corresponding referring bounding boxes.

## B    Details of CODE Benchmark

In this section, we provide comprehensive details about how we collect the CODE dataset to facilitate research on contextual object detection.

**Data Format.** Our CODE benchmark follows the data format of the COCO dataset and includes additional fields to facilitate the evaluation, as shown in Figure 5. The images and annotations used in our new benchmark are based on Flickr 30k [78] and Flickr30k Entities [53]. We tokenize the language caption using the LLM tokenizer and record the related language tokens. For each object name that appears in the tokens generated by the tokenizer, we track the start and end indices, which will be replaced with the [MASK] token for our contextual cloze test task.

**Word Clouds.** In the contextual cloze test setting, our CODE dataset consists of 10,346 unique object words that are masked and required to be predicted. Figure 6 presents the word cloud visualizations of object words in our dataset. We can observe both high-frequency words such as 'man' and 'woman,' as well as low-frequency words such as 'player', 'scooty', and 'breadstick,' which pose challenges for accurate predictions. Therefore, achieving precise predictions for these object words requires understanding contextual information.

## C    Details of Evaluation for Contextual Cloze Test

Existing object detection datasets, such as Pascal VOC [18], Microsoft COCO [38], Open Images [34], LVIS [22], Objects365 [62] and V3Det [69], rely on predefined mappings between label IDs and class names for evaluation purposes. For example, the COCO dataset uses a mapping like (1, person), (2, bicycle), . . . , (80, toothbrush) for its 80 classes. As shown in Fig.7(a), in order to be classified as true positives, predicted bounding boxes must exhibit both high IoU overlap and identical class IDs

Table 5: Comparison of our proposed three contextual object detection settings with previous related tasks.

| Tasks | Language Input | Output(s) | Remark |
|---|---|---|---|
| Object Detection | ✗ | box, class label | pre-defined class labels |
| Open-Vocabulary Object Detection | (optional) class names for CLIP | box, class label | pre-defined class labels |
| Referring Expression Comprehension | complete referring expression | box that expression refers to | / |
| **Contextual Cloze Test** (ours) | **incomplete** expression object names are masked | {box, **name**} to complete the mask | **name** could be most valid English word |
| Image Captioning | ✗ | language caption | |
| **Contextual Captioning** (ours) | ✗ | language caption, **box** | |
| Visual Question Answering | language question | language answer | |
| **Contextual QA** (ours) | language question | language answer, **box** | |



(a) Data format

(b) Exemplar sample

(c) We save the start and end indices of tokens that represent the object name in the "token_ids" field of annotations

Figure 5: Our CODE benchmark follows the data format of the COCO dataset [38], with additional fields (blue color) including the language caption, token ids, and object name. Token ids record the start and end position index of the object name existing in the language tokens.

to the ground-truth boxes. In certain scenarios, such as zero-shot[4] or open-vocabulary [82] object detection settings, the predefined classes are divided into two separate groups: *base* and *novel*, to evaluate the model's generalization capability. However, these evaluations still rely on the predefined ID-name mappings, while objects with names not included in predefined mappings are impossible.
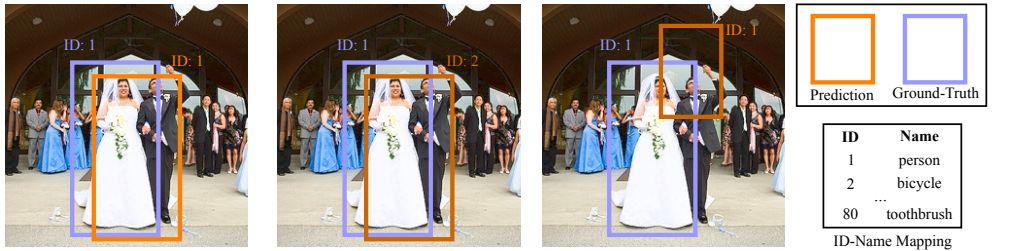
Human perception does not depend on pre-defined class IDs. Therefore, for our proposed contextual cloze test task, we have established new evaluation criteria that use object names from human language. In this evaluation, given a masked language expression and the indexes of the masked words, we classify predicted boxes as true positives if they i) exhibit high IoU overlap, ii) share the same meaning, and iii) have an identical masked index as the ground truth boxes. Conversely, predictions are considered false positives. The masked indexes are employed to differentiate cases where multiple objects have the same name but are located at different [MASK] token positions within a sentence. The object names correspond to the most valid English words decoded by the Tokenizer of LLMs.

After defining our name-based criteria as true-positive/false-positive metrics, we could compute the overall Average Precision (AP) metric for evaluation. We follow the COCO dataset to set the IoU thresholds ranging from 0.5 to 0.95 with a step size of 0.05. The per-name AP is not computed because there are numerous long-tailed infrequent names, of which only a few examples are available for evaluation.
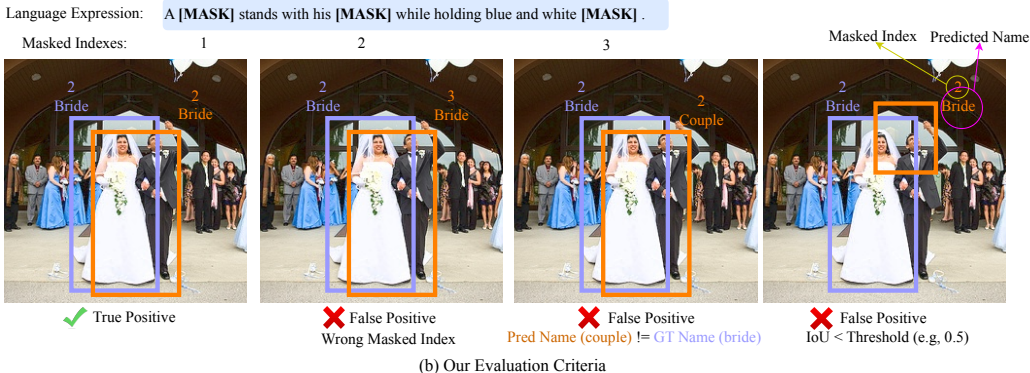
**AP@5 for Top-5 Predicted Names.** In some cases, our evaluation metric can be overly stringent, particularly when dealing with numerous synonyms or fine-grained categories that are challenging for annotators to distinguish. Similar challenges have been encountered in previous image classification datasets like ImageNet [12], where the top-5 accuracy metric is used as a supplementary metric to the top-1 accuracy metric. Therefore, we also introduce a supplementary metric called top-5 AP (AP@5), which relaxes the definition of true positives. Under AP@5, if the ground-truth name is among the top-5 predictions, the predictions are considered true positives. In contrast, the AP metric calculated based on the top-1 prediction result is referred to as AP@1 to differentiate it from AP@5.

Figure 6: Word clouds of object words presented in the CODE train set (**left**) and test set (**middle**, **right**). The **middle** figure represents the visualization of high-frequency words in the test set, while the **right** figure showcases the visualization of low-frequency words.



(a) Evaluation Criteria for Previous Detection Datasets

(b) Our Evaluation Criteria

Figure 7: The comparison of (**a**) the evaluation criteria for the traditional object detection task, and (**b**) evaluation of our contextual cloze test.

**Implementation Details.** We modify the famous *pycocotools* package [2] provided in the COCO dataset and create the evaluation script.

# D  More Experiments

In this section, we first demonstrate that our ContextDET could be extended to segmentation tasks, such as the referring image segmentation task discussed in Section D.1. Next, we provide additional qualitative results, including failure cases, in Section D.2.

## D.1  Referring Image Segmentation

Our ContextDET is not limited to object detection and can be extended to the image segmentation task, in which the goal is to assign a pixel-level label to each pixel in the input image. To adapt our ContextDET framework for segmentation, we introduce an extra pixel-level segmentation head that takes the full visual tokens $c$ as inputs. To train the segmentation model, we use a pixel-wise cross-entropy loss $\mathcal{L}_{\text{mask}}$ and Dice loss $\mathcal{L}_{\text{dice}}$, where ground-truth labels are pixel-level masks for matched objects in an image.

We choose the referring image segmentation task as a representative benchmark to evaluate the segmentation performance of ContextDET. The referring image segmentation task aims to segment

---

[2]https://github.com/cocodataset/cocoapi

Table 6: Comparisons with state-of-the-art methods on three referring image segmentation benchmarks in terms of the mean Intersection over Union (mIoU) metric.

| # | Method | Venue | Language Model | Vision Backbone | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|--------|-------|----------------|-----------------|---------|---|---|----------|---|---|----------|---|
| | | | | | val | testA | testB | val | testA | testB | val | test |
| 1 | VLT [13] | ICCV'21 | Bi-GRU | DN53 | 65.65 | 68.29 | 62.73 | 55.50 | 59.20 | 49.36 | 52.99 | 56.65 |
| 2 | SeqTR [88] | ECCV'22 | Bi-GRU | DN53 | 71.70 | 73.31 | 69.82 | 63.04 | 66.73 | 58.97 | 64.69 | 65.74 |
| 3 | RefTR [48] | NeurIPS'21 | BERT-base | RN101 | 74.34 | 76.77 | 70.87 | 66.75 | 70.58 | 59.40 | 66.63 | 67.39 |
| 4 | LAVT [76] | CVPR'22 | BERT-base | Swin-B | 74.46 | 76.89 | 70.94 | 65.81 | 70.97 | 59.23 | 63.34 | 63.62 |
| 5 | PolyFormer [40] | CVPR'23 | BERT-base | Swin-B | 75.96 | 77.09 | 73.22 | 70.65 | 74.51 | 64.64 | 69.36 | 69.88 |
| 6 | ContextDET | - | OPT-2.7B | Swin-B | **76.40** | **77.39** | **74.16** | **71.67** | **75.14** | **65.52** | **69.89** | **70.33** |

regions described by fine-grained input language query. Language queries will act as conditional inputs for the visual decoder in ContextDET. We use three commonly-used datasets: RefCOCO [79], RefCOCO+ [79] and RefCOCOg [49]. On RefCOCO and RefCOCO+, we follow the default training/validation/testA/testB data split in Yu *et al* [79]. For RefCOCOg, we use the RefCOCO-umd splits [49]. We report the mean Intersection over Union (mIoU), which is calculated by averaging the IoU scores across all test samples.

We compare ContextDET with some state-of-the-art methods in Table 6. ContextDET achieves better results with mIoU gains of 0.63% and 0.45% on the validation/test splits over PolyFormer [40].

### D.2  More Qualitative Results

We provide more qualitative results predicted by ContextDET in the contextual cloze test (Figure 8), contextual captioning (Figure 9), and contextual QA settings (Figure 10). The selected images are sourced randomly from the web and are not included in the training data. We observe that ContextDET effectively predicts contextual object words, including terms like 'teacher', 'student', 'doctor', and 'nurse', along with their corresponding bounding boxes. In addition, we find some failure cases. For instance, the predicted object words may be incorrect, particularly for less common terms like 'earth'. Our ContextDET is less robust when it comes to occluded objects, such as 'sheep'. We aim to address these limitations in future research.

## E  Broader Impact

In this paper, we propose ContextDET that is capable of detecting objects within multimodal vision-language contexts. Our proposed ContextDET could facilitate the development of more real-world applications with human-AI interaction, such as AR/VR systems and robotics. However, relying on LLMs in our method may raise potential concerns when LLMs are not used properly. For example, LLMs could be used for harmful applications, such as plagiarism, spamming, and telemarketing fraud. Therefore, further research focused on detecting AI-generated text [45, 61] is essential to address these issues and mitigate the negative impacts.

Figure 8: Qualitative examples of the **contextual cloze test**.



Figure 9: Qualitative examples of the **contextual captioning**.

Figure 10: Qualitative examples of the **contextual QA**.