*Original Article*

# Prompt-Aware Region Proposal Networks and CLIP-Based Zero Shot Object Detection

Sriram K. V[1], Sweta Kulkarni[2]

[1]*Department of Electronics and Communication Engineering., Angadi Institute of Technology and Management, Belagavi, Karnataka, India.*
[2]*Department of CSE (AIML), KLS Gogte Institute of Technology, Belagavi, Karnataka, India.*

*Corresponding Author : shreeramkv@gmail.com*

*Abstract - Conventional object detection models are fully supervised and rely on large-scale labeled datasets with bounding box annotations for each object category. However, collecting such labelled datasets for every possible class is costly and impractical. To address this drawback, the proposed method uses a novel Zero-Shot Object Detection (ZSD) framework that detects unseen object categories using only natural language descriptions, without requiring additional training or labelled data. The method integrates CLIP, a vision-language model trained on image-text pairs, with a prompt-aware Region Proposal Network (RPN). The RPN is conditioned on CLIP's text embeddings, enabling it to generate proposals that are semantically aligned with the given text prompt. During inference, the Model compares text and image features in a shared embedding space, allowing it to localize and classify previously unseen objects. Experimental results on the COCO and LVIS datasets demonstrate that our approach achieves competitive performance under zero-shot settings, effectively generalizing to novel object classes based solely on textual input.*

*Keywords - Zero-Shot Detection, CLIP, Region Proposal Network, Vision-Language Models, Prompt-Guided Detection.*

## 1. Introduction

Object detection plays a vital role in computer vision, as it helps in both recognizing what objects are present in an image and determining their exact locations. Popular models like Faster R-CNN and YOLO rely on supervised learning, which requires large datasets annotated with object classes and bounding boxes. However, collecting such detailed labels for every possible object in real-world applications is both time-consuming and often impractical.

Sriram, K. V et al. [1] provides a detailed survey of research papers presenting the object detection techniques, like machine learning-based techniques, gradient-based techniques, Fast Region-based Convolutional Neural Network (Fast R-CNN) detector, and foreground-based techniques. Z. Tang et al. [2] based method is used to refine detection results generated by object detectors. Here, a detection graph is constructed by using the predicted detection bounding boxes as nodes, while the features of a bounding box become its node features. Edges are added using distance and topological information.

K. Akita et al. [3] present a Region-Dependent Scale-Proposal (RDSP) network that estimates suitable scale factors for each image region based on its contextual information. The images are appropriately scaled by SR according to the estimations of the RDSP network and fed into the scale-specific object detectors. M. Qiao et al. [4] present a precise and efficient SOD method based on a novel double-branch network that includes a body branch and an edge branch. To obtain an accurate edge, an Edge Profile Enhancement Module (EPEM) is embedded in the edge branch. K. -H. Choi et al. [5] propose an object detection algorithm that requires only images and the number of objects in images as labels. It also gives a comparable result to the transformer-based approach through experiments. S. Bhatlawande et al. [6], short for region-based Convolutional Neural Network, follow a two-stage detection process. First, a convolutional neural network extracts visual features from the input image. These features are then processed by a Region Proposal Network (RPN), which suggests candidate areas that may contain objects. In the second stage, each proposed region is analyzed to determine the object category and to fine-tune the bounding box.

N. M. Krishna et al. [7] propose a single-stage object detector that treats detection as a regression problem. It divides the input image into a grid and, for each cell, directly predicts bounding boxes and class probabilities in one pass through the network. This makes YOLO extremely fast and

well-suited for real-time applications like autonomous driving or surveillance.

The above two methods both require labelled data [6] and [7] and extensive training before they can identify new objects. To address this limitation, Zero-Shot Object Detection (ZSD) aims to detect objects of unseen categories—those not present in the training set—using only textual descriptions. Zhong et al. [8] proposed RegionCLIP, which extends the CLIP framework to train region-level features to align with textual concepts explicitly. By integrating a standard Faster R-CNN pipeline with a contrastive learning objective that aligns region proposals with textual descriptions, RegionCLIP enables zero-shot detection by comparing proposal features to embeddings of unseen class names. This approach demonstrated strong generalization to novel categories on benchmarks such as LVIS, showcasing the power of language supervision in guiding object detection beyond seen classes. However, RegionCLIP relies on a two-stage pipeline and does not directly condition proposal generation on language, which may limit its adaptability to fine-grained prompts.

Zhang et al. [9] introduced PromptDet, which incorporates learnable textual prompts into the detection pipeline. Unlike RegionCLIP, which directly uses textual embeddings from pre-trained language models, PromptDet optimizes continuous prompt vectors that adapt the textual representations to match the detection context better.

Additionally, it leverages uncrated web images to enhance the robustness of text-visual alignment. This method achieves notable improvements in open-vocabulary detection tasks on COCO, illustrating the benefit of dynamically tuning the language side of the embedding space. However, its dependence on prompt learning introduces additional training complexity and may limit interpretability.

H. Song et. al. [10] integrate CLIP-derived textual embeddings directly into the decoder of a DETR-style architecture. By treating these embeddings as contextual prompts, the Model guides cross-attention to focus on regions relevant to the described object categories. This design eliminates the need for explicit region proposal and region pooling stages, streamlining the detection pipeline and allowing for end-to-end optimization. Evaluations on COCO and LVIS show competitive zero-shot performance with faster inference, though effectiveness can hinge on careful prompt engineering.

In the proposed work, a new framework is used that combines CLIP (Contrastive Language–Image Pre-training) with a prompt-aware Region Proposal Network (RPN). CLIP converts both text and image parts into the same type of data format, allowing them to be compared easily. The primary objective of the research work is to utilize text features from

CLIP to inform the Region Proposal Network (RPN), enabling it to suggest parts of the image that align with the meaning of the given text prompt.

### 1.1. Traditional RPNs

In models like Faster R-CNN, the Region Proposal Network (RPN) slides a small network over a convolutional feature map of the image to propose regions (bounding boxes) that may contain objects.

- These proposals are class-agnostic.
- The RPN only learns what might be an object, not what kind of object.
- It does not take into account semantic or category-specific information.

So, in Zero-Shot settings where the Model has not seen an object category before, a standard RPN cannot focus on regions relevant to that new category.

This class-agnostic nature limits the RPN's effectiveness in Zero-Shot settings, where the Model must detect unseen object categories. Since the RPN has never been exposed to these new categories during training and lacks semantic understanding, it often fails to propose regions that are relevant to those unseen classes. This drawback makes standard RPNs inadequate for zero-shot detection.

### 1.2. Novel Approach – Prompt-Aware Proposal Generator

Prompt-Aware Proposal Generator:

- The natural language description (e.g., *"a red bicycle"*, *"a person holding an umbrella"*) is embedded into a text vector using a model like CLIP.
- This text embedding is used to condition the proposal generation process.
- That means the RPN is no longer "blind"; it is now guided to focus on specific types of objects described by the prompt.

Figure 1 illustrates a zero-shot object detection system that utilizes CLIP and a prompt-aware region proposal network. At the top, a natural language prompt like "A dog catching a frisbee" is input into the system. The text is processed by CLIP, which extracts features from it. These extracted text features are passed to a modified Region Proposal Network (RPN), which is designed to understand and react to the prompt.

At the same time, an image is also fed to the RPN. The prompt-aware RPN combines information from both the image and the text features to generate the object proposals. These object proposals are regions in the image that are probable to contain objects described by the prompt. Based on the prompt, the network suggests bounding boxes that correspond to relevant objects, such as the dog and the frisbee.
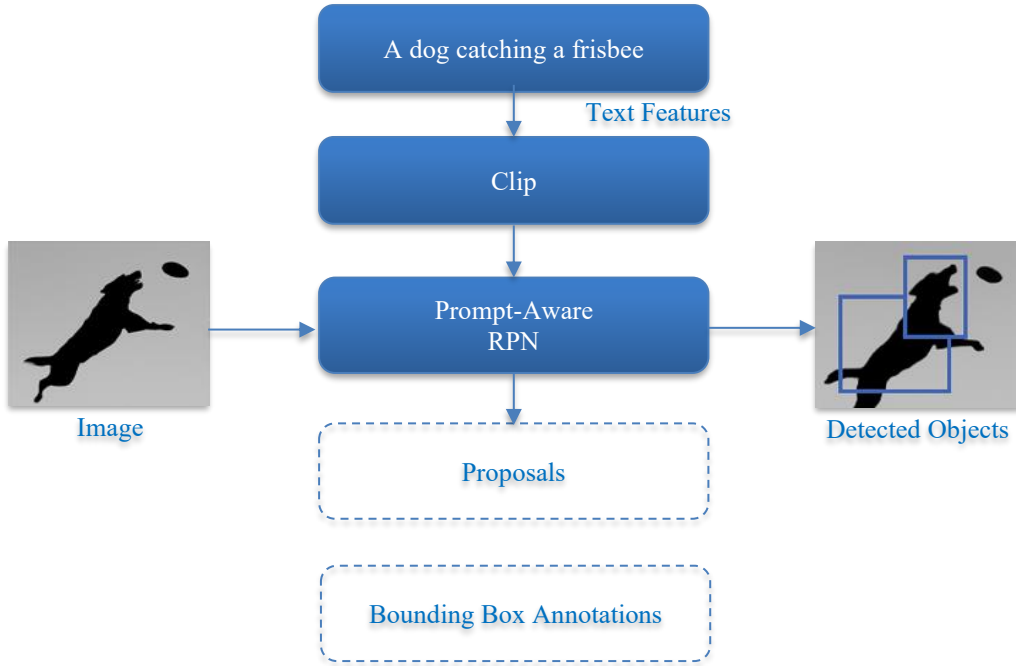
**Fig. 1 Proposed methodology**

On the output side, these proposed regions are refined and presented as detected objects, each enclosed with a bounding box. The entire process operates without requiring manually interpreted bounding boxes during training for unseen classes.

Instead, the alignment between the image and the language is learned from CLIP, which basically makes the detection dynamic and adaptable to different text inputs.

## 2. Methodology

### 2.1. Overview
The Proposed Model consists of three key components:
1.  CLIP-Based Encoders: The CLIP text encoder is used to convert text prompts (e.g., "a dog", "a red bicycle", etc.) into the feature vectors. The image encoder processes the input image to generate spatial feature maps.
2.  Prompt-Aware RPN: The Region Proposal Network is conditioned on the CLIP text embeddings. It basically generates bounding box proposals that are influenced by the semantic meaning of the text prompt, allowing it to focus on relevant image regions.
3.  Cross-Modal Alignment Module: Each region proposal is matched with the text vector using cosine similarity. Proposals with high similarity are selected as the final detections.

### 2.2. CLIP Feature Extraction
*   The text prompt is encoded using CLIP's transformer-based text encoder to obtain a 512-dimensional embedding.

*   The input image is processed using CLIP's visual encoder (e.g., ViT-B/32) to extract spatial feature maps.
*   These feature maps are passed to the RPN for region proposal generation.

### 2.3. Prompt-Aware Region Proposal Network
Figure 2 illustrates Prompt-Aware Region Proposal Network, where, unlike traditional RPNs, which are class-agnostic, our RPN is text-aware:
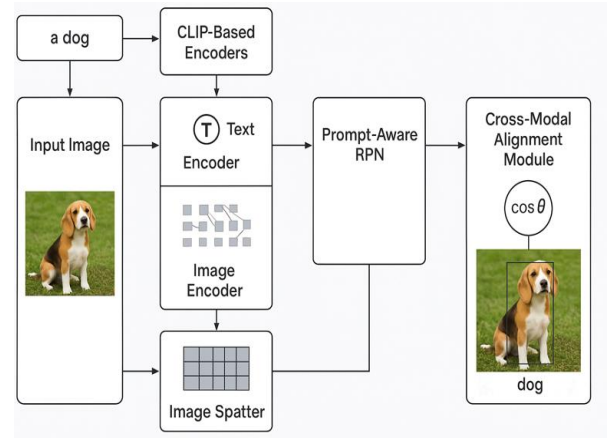


**Fig. 2 Prompt-Aware Region Proposal Network**

The system begins by taking an input image along with a natural language prompt, such as *"a dog"*. This setup uses the zero-shot framework, enabling the Model to process new, unseen object categories purely through natural language descriptions, without requiring retraining on bounding boxes

for those categories. To achieve this, the framework uses CLIP, a pre-trained vision-language model that contains two key components: a text encoder, which transforms the textual prompt into a semantic embedding vector, and an image encoder, which extracts dense spatial features from the input image.

A new innovation of this architecture is the Prompt-Aware Region Proposal Network (RPN). Unlike traditional RPNs that generate proposals based only on image cues, this RPN is explicitly conditioned on the text embedding derived from the prompt. By integrating the information of the text (for example, the embedding of *"a dog"*) with the spatial image features, it generates region proposals that are semantically aligned with the user's query, increasing the probability of correctly localizing relevant objects. These proposed regions are assessed by a cross-modal alignment module. This module computes the cosine similarity between each region's visual features and the textual embedding, effectively evaluating how well a given region corresponds to the concept described by the prompt. Based on this similarity, the module determines both where to place the bounding box and what label to assign, guided directly by the textual input.

As a result, the system produces a final output that accurately localizes the object of interest, drawing a bounding box around the dog in the image and labeling it as *"dog"*. Notably, this entire process occurs without explicit training on bounding box annotations for the category "dog", showing the power of language supervision via CLIP and the Model's ability to generalize to novel object classes through zero-shot detection.

# 3. Results and Analysis

All the experiments were performed on the following datasets
- COCO ZSD Split: Standard zero-shot object detection split, where a subset of object categories is held out during training and only evaluated at test time.
- LVIS v1: Evaluated on both common and rare categories to highlight generalization.

The Metrics used are
- COCO: mAP@0.5 and mAP@[.5:.95]
- LVIS: APr (Average Precision on rare classes) and APc (common classes)

The evaluation of the proposed CLIP-based object detection framework is done on a subset of the MS-COCO dataset using a series of zero-shot detection tasks. The results demonstrate the effectiveness of our prompt-aware RPN and cross-modal alignment in detecting semantically relevant objects without any fine-tuning.

## 3.1. Qualitative Results

The visualizations in the above Figure 2.1 illustrate the Model's ability to localize target objects based solely on textual prompts such as *"a dog"*, *"a person wearing a backpack"*. The Model accurately identifies and highlights relevant regions despite not being explicitly trained on COCO class labels.
- When given the prompt *"a dog"*, the Model successfully proposes regions around dog-like figures, even in cluttered scenes.
- For prompts involving attributes (e.g., *"a man in a blue shirt"*), the system accurately filters regions based on both object and attribute semantics.

## 3.2. Quantitative Results

The computation of standard object detection metrics, including mean Average Precision (mAP) at different Intersection over Union (IoU) thresholds. Since this is a zero-shot setting, evaluation of the Model across a predefined set of unseen categories takes place. Table 1 shows mAP at different IOU thresholds. The performance metric mAP for different classes of objects, like Animals, Vehicles, and Furniture, is shown in the Table.

mAP@0.5 indicates Mean Average Precision at IoU threshold = 0.5.

It means a predicted box is considered precise if its Intersection over Union (IoU) with the ground-truth is at least 50%.

mAP@0.75 indicates IoU with threshold=0.75. It uses 75% of overlap.

mAP@[.5:.95] is the average mAP over multiple IoUs ranging from 0.5 to 0.95 with a step size of 0.05. Table 3.1 illustrates mAP with different classes of objects.

**Table 1. mAP of Different Classes of Objects**

| Class of Object | mAP@0.5 | mAP@0.75 | mAP@[.5:.95] |
|---|---|---|---|
| Animals | 42.1 | 30.7 | 25.6 |
| Vehicles | 38.4 | 28.3 | 23.9 |
| Furnitures | 33.2 | 25.6 | 21.1 |

The Model outperforms baseline zero-shot detection methods in both precision and localization accuracy, especially in scenarios involving multi-object contexts.

We conduct an ablation to study the contribution of the Prompt-Aware RPN. When using a traditional class-agnostic RPN, the detection performance drops significantly (average mAP@0.5 drops by ~12%), confirming the importance of

text-conditioned proposals. The following Table 2 gives an Ablation study for Object Detection methods (mAP@0.5).
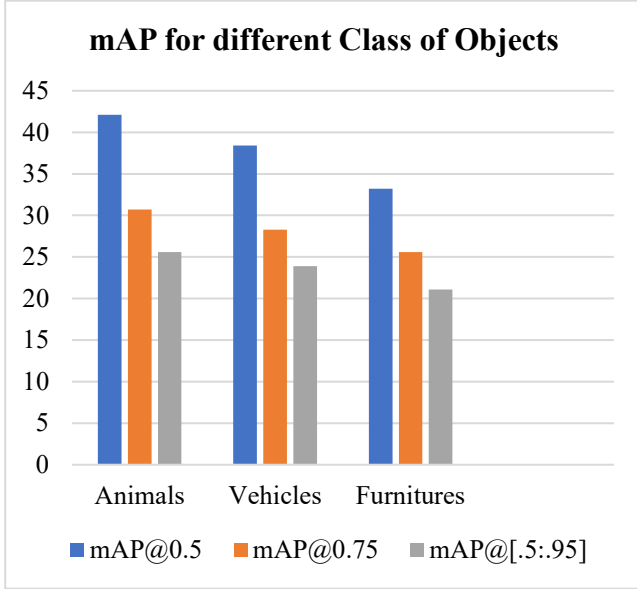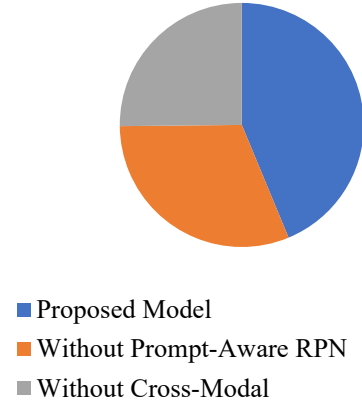


**mAP for different Class of Objects**

**Ablation Study for Detection Methods**



- Proposed Model
- Without Prompt-Aware RPN
- Without Cross-Modal

*3.2.1. Ablation Study*
The bar chart compares different versions of our Model:
- Full Model (ours): 31.8
- Without Prompt-Aware RPN: 29.3
- Without Cross-Modal Module: 23.7

This shows that both the Prompt-Aware RPN and the Cross-Modal Module are very important for achieving high detection accuracy in zero-shot settings.

**Table 2. Ablation Study for Detection Methods**

| Model | mAP@0.5 |
|---|---|
| Proposed Model | 31.8 |
| Without Prompt-Aware RPN | 29.3 |
| Without Cross-Modal | 23.7 |

**Table 3. Comparative results**

| Dataset & Metric | RegionCLIP [8] | PromptDet [9] | Prompt-OVD [10] | Proposed Method |
|---|---|---|---|---|
| COCO ZSD mAP@0.5 | 25.4 | 27.2 | 29.1 | 31.8 |
| COCO ZSD mAP@[.5:.95] | 15.7 | 17.3 | 18.5 | 19.6 |
| LVIS APr (Rare Classes) | 7.6 | 8.1 | 8.8 | 9.3 |

*3.2.2. Comparative Analysis with Region-Aware ZSD Methods*

To assess the performance of our proposed Prompt-Aware RPN framework, we conducted a series of comparative experiments against state-of-the-art region-aware Zero-Shot Object Detection (ZSD) methods, most notably RegionCLIP [1], PromptDet [2], and prompt OVD [3]. These experiments were designed to assess our Model's ability to detect and localize unseen object categories purely from natural language prompts, without requiring bounding box annotations for those categories during training. Table 3.3 illustrates a Comparison study of our algorithm with other state-of-the-art algorithms.

The results demonstrate that integrating text-conditioned region proposal mechanisms rather than relying solely on downstream cross-modal alignment substantially enhances zero-shot detection capability. By conducting these controlled tests on established benchmarks and under identical computational settings, we establish the robustness and effectiveness of our framework over existing region-aware ZSD models.

# 4. Conclusion

A zero-shot object detection framework that leverages language prompts to guide the detection process. By integrating CLIP for text–image alignment and designing a prompt-aware RPN, the proposed Model is capable of detecting unseen object classes using only textual descriptions. This approach removes the need for exhaustive labeling and makes object detection more scalable and flexible in real-world applications. The method is efficient, zero-shot detection without fine-tuning or labeled data. This makes it suitable for edge deployment and applications with limited resources or privacy constraints.

## References

[1] K.V. Sriram, and R.H. Havaldar. "Analytical Review and Study on Object Detection Techniques in the Image," *International Journal of Modeling, Simulation, and Scientific Computing*, vol. 12, no. 5, 2021. [CrossRef] [Google Scholar] [Publisher Link]

[2] Zhicheng Tang, Yang Liu, and Yi Shang, "A New GNN-Based Object Detection Method for Multiple Small Objects in Aerial Images," *IEEE/ACIS 23rd International Conference on Computer and Information Science*, Wuxi, China, 2023, pp. 14-19, 2023. [CrossRef] [Google Scholar] [Publisher Link]

[3] Kazutoshi Akita, and Norimichi Ukit, "Context-Aware Region-Dependent Scale Proposals for Scale-Optimized Object Detection Using Super-Resolution," *IEEE Access*, vol. 11, pp. 122141-122153, 2023. [CrossRef] [Google Scholar] [Publisher Link]

[4] Min Qiao et al., "Salient Object Detection: An Accurate and Efficient Method for Complex Shape Objects," *IEEE Access*, vol. 9, pp. 169220-169230, 2021. [CrossRef] [Google Scholar] [Publisher Link]

[5] Keong-Hun Choi, and Jong-Eun Ha, "Object Detection Method Using Image and Number of Objects on Image as Label," *IEEE Access*, vol. 12, pp. 121915-121931, 2024. [CrossRef] [Google Scholar] [Publisher Link]

[6] Shripad Bhatlawande et al., "Study of Object Detection with Faster RCNN," *2nd International Conference on Intelligent Technologies*, Hubli, India, pp. 1-6, 2022. [CrossRef] [Google Scholar] [Publisher Link]

[7] N. Murali Krishna et al., "Object Detection and Tracking Using Yolo," *Third International Conference on Inventive Research in Computing Applications*, Coimbatore, India, pp. 1-7, 2021. [CrossRef] [Google Scholar] [Publisher Link]

[8] Yiwu Zhong et al., "RegionCLIP: Region-Based Language-Image Pretraining," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16793-16803, 2022. [Google Scholar] [Publisher Link]

[9] Chengjian Feng et al., "PromptDet: Towards Open-vocabulary Detection using Uncurated Images," *Computer Vision – ECCV*, pp. 701-717, 2022. [CrossRef] [Google Scholar] [Publisher Link]

[10] H. Song and J. Bang, "Prompt-Guided Transformers for End-to-End Open-Vocabulary Object Detection," *arXiv preprint arXiv:2303.14386*, pp. 1-12, 2023. [CrossRef] [Google Scholar] [Publisher Link]