

# Learning Cross-Modal Affinity for Referring Video Object Segmentation Targeting Limited Samples

Guanghui Li<sup>1\*</sup>, Mingqi Gao<sup>2,3\*</sup>, Heng Liu<sup>1†</sup>, Xiantong Zhen<sup>4</sup>, Feng Zheng<sup>2†</sup>  
<sup>1</sup> Anhui University of Technology, <sup>2</sup> Southern University of Science and Technology,  
<sup>3</sup> University of Warwick, <sup>4</sup> United Imaging  
 guanghui.li1998@gmail.com, mingqi.gao@outlook.com, hengliusky@aliyun.com,  
 zhenxt@gmail.com, f.zheng@ieee.org

## Abstract

Referring video object segmentation (RVOS), as a supervised learning task, relies on sufficient annotated data for a given scene. However, in more realistic scenarios, only minimal annotations are available for a new scene, which poses significant challenges to existing RVOS methods. With this in mind, we propose a simple yet effective model with a newly designed cross-modal affinity (CMA) module based on a Transformer architecture. The CMA module builds multimodal affinity with a few samples, thus quickly learning new semantic information, and enabling the model to adapt to different scenarios. Since the proposed method targets limited samples for new scenes, we generalize the problem as *- few-shot referring video object segmentation (FS-RVOS)*. To foster research in this direction, we build up a new FS-RVOS benchmark based on currently available datasets. The benchmark covers a wide range and includes multiple situations, which can maximally simulate real-world scenarios. Extensive experiments show that our model adapts well to different scenarios with only a few samples, reaching state-of-the-art performance on the benchmark. On Mini-Ref-YouTube-VOS, our model achieves an average performance of 53.1  $\mathcal{J}$  and 54.8  $\mathcal{F}$ , which are 10% better than the baselines. Furthermore, we show impressive results of 77.7  $\mathcal{J}$  and 74.8  $\mathcal{F}$  on Mini-Ref-SAIL-VOS, which are significantly better than the baselines. Code is publicly available at [https://github.com/hengliusky/Few\\_shot\\_RVOS](https://github.com/hengliusky/Few_shot_RVOS).

## 1. Introduction

Referring video object segmentation (RVOS) aims to segment target objects described in natural language in

\*Equal contribution. This work was done when G. Li visited to Feng Zheng Lab in Southern University of Science and Technology.

†Corresponding author.

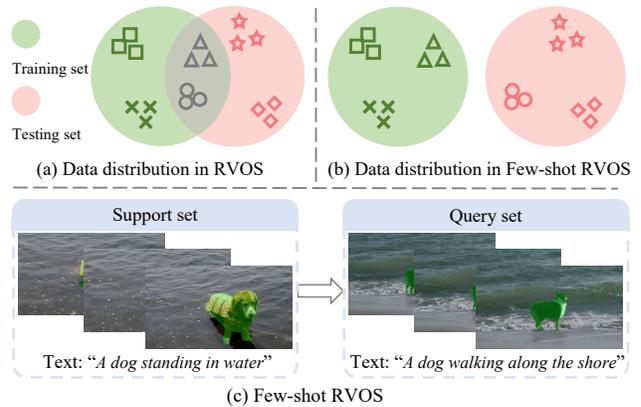


Figure 1: Comparison of Few-shot RVOS and RVOS and the setting of Few-shot RVOS. (a) The training and testing sets overlap in the RVOS. (b) Disjoint training and testing sets in the Few-shot RVOS. Different shapes represent different classes. (c) Few-shot RVOS segments the referred object of the same class as the support set in the video.

videos. In real-world scenarios, it has a wide range of applications, such as video editing [9] and human-computer interaction, so RVOS has attracted much attention from the research community. Unlike traditional semi-supervised video object segmentation [10], RVOS is more challenging because it not only lacks the ground-truth mask of the first frame of the video but also needs to interact with multimodal information of vision and language.

The great success of various tasks based on deep learning benefits from sufficient labeled data. Detailed annotated masks and language descriptions in real-world RVOS tasks are relatively scarce. The researchers have to annotate each frame in the video in detail and provide a referring expression for the segmentation object. Therefore, obtaining high-quality labeled data requires a high cost. With the popularity of movies, YouTube videos, TikTok streaming videos,

etc., video data in various fields has shown explosive growth in this media age. The demand for processing diverse data has brought significant challenges. In order to handle diverse data, existing RVOS methods must rely on massive and diverse labeled data for training. But due to fixed and limited training classes, existing RVOS methods [1, 33] are essentially constrained to adapt to the highly dynamic and highly diverse data in the real world. On the other hand, if ones fine-tune existing RVOS methods on a few samples to adapt to real-world data, high-quality results are hard to be achieved because the labelled data is insufficient to support the model for learning the new semantics. Therefore, how to make the RVOS methods applicable to real-world diverse data with a lower cost is an urgent problem.

To address this problem, we propose the cross-modal affinity (CMA) module to build multimodal relationships in a few samples and learn new semantic information for diversified data. Specifically, given only a few annotated samples (language expressions and the referred object masks), we hierarchically fuse visual and text features in a cross-attention manner to obtain robust feature representations for a specific category. In this way, the model can handle enormous data in the same category more efficiently.

Essentially, the proposed method targets limited samples. Therefore, we generalize the problem as Few-Shot Referring Video Object Segmentation (FS-RVOS). We show the setting of FS-RVOS and the difference from existing RVOS in Figure 1. Unlike RVOS, the training and testing sets' categories disjoint in the FS-RVOS. Given a few support video clips together with corresponding language descriptions and object masks, FS-RVOS aims at segmenting videos in the query set, as shown in Figure 1(c).

The key to FS-RVOS lies in the support set utilization and understanding of vision-language information. To better leverage the information in the support set, two methods have been proposed based on the prototype and attention mechanisms. The prototype-based [21, 35] methods compress the features belonging to different classes to obtain prototypes. However, noise is easily generated during the process. In addition, the spatial structures are ignored, resulting in different degree of information loss. Another methods [37, 38, 39] employ the attention mechanism to encode foreground pixels from support features and aggregate them with query features. Although these methods achieve high-quality results in image and video domains, they are still under-explored in vision-language tasks.

To better utilize vision-language inputs, we propose the cross-modal affinity module to build the multimodal relationships between samples in the support and query sets. Specifically, multimodal features within the support set and query set are first fused separately. The information among them is then aggregated, which effectively prevents query features from being biased by irrelevant features.

Since this is the first work exploring Few-shot RVOS, the existing datasets are not directly applicable. Therefore, we build up a new FS-RVOS benchmark based on Ref-YouTube-VOS [22], named Mini-Ref-YouTube-VOS. The new benchmark covers a wide range with a balanced number of high-quality videos in each category. To measure the model's generalization ability, we also build a dataset different from natural scenes based on a synthetic dataset SAIL-VOS [13], named Mini-Ref-SAIL-VOS. Since only videos and detailed annotated masks exist in the SAIL-VOS dataset, we add natural language descriptions corresponding to the segmentation targets for the dataset.

The main contributions of this work are as follows.

- For real-world limited samples, we propose a **Cross-Modal Affinity (CMA)** for building multimodal information affinity for referring video object segmentation.
- We explore a novel Few-shot RVOS problem, which learns new semantic information with limited samples and can adapt to diverse scenarios.
- We build up the first FS-RVOS benchmark, where we conduct comprehensive comparisons with existing methods, showing the superiority of the proposed model.

## 2. Related Work

**Few-Shot Semantic Segmentation.** Few-shot semantic segmentation, first proposed by Shaban et al. [23], aims to learn how to segment new categories of images through a few samples. Recent advances in few-shot semantic segmentation originate from the application of metric learning. Based on PrototypicalNet, Dong et al. [7] first employ the metric learning technique and apply cosine similarity between pixels and prototypes for prediction. In addition, PANet et al. [31] introduce prototype alignment regularization to simplify the framework. PFENet et al. [28] use prior knowledge from the pretrained backbone to find the regions of interest and the different designs of feature pyramid modules, and previous leverage mappings to achieve better segmentation performance.

However, the effectiveness of these few-shot segmentation methods depends mainly on the quality of the prototypes obtained from the support set. Fan et al. [8] address the critical intra-class appearance differences inherent in the few-shot segmentation problem by performing self-support matching with query features. Their strategy effectively captures the consistent underlying features of query objects to match query features. Tian et al. [27] propose a novel context-aware prototype learning method that leverages prior knowledge from support samples and dynamically enriches contextual information by using adaptive features. Motivated by the idea that using the base learner to identify confusing regions in the query image and further

refining the predictions of the meta-learner, BAM [15] establishes a new method for few-shot segmentation that does not focus on feature extraction or visual correspondence.

Compared to image-based few-shot segmentation, few-shot video object segmentation works are relatively rare and remain in the early stage. The initial work [25], [3] mainly solves this problem through the attention mechanism. However, these methods do not consider temporal information. Thus, based on temporal transductive, the recent work [24] applies reasoning mechanisms and has achieved good results in cross-domain scenarios. In a word, all the above-mentioned few-shot segmentation methods are only for a single modality - image or video, and cannot handle segmentation under multimodal conditions (i.e., with linguistic referring expressions).

**Referring Video Object Segmentation.** Gavriluk et al. [11] first introduce the RVOS task. They generate convolutional dynamic filters from textual representations and convolve them with visual features of different resolutions to obtain segmentation masks. To overcome the limitations of traditional dynamic convolution, Wang et al. [30] propose a context-modulated dynamic convolution operation for RVOS, where the kernel is generated from language sentences and surrounding contextual features. However, since the focus is only on video actors and actions, their approach only applies to a few object classes and action-oriented descriptions. Weak-Shot Semantic Segmentation (WSSS) [4, 42] focuses on the overall scene in the image, treating masks and text as the support set and text as the query set. However, in WSSS, the text is limited to single words or phrases indicating class names and directly mapped to labels for pixel-level classification.

Khoreva et al. [14] propose a two-stage approach that first performs referring expressions grounding and then utilizes the predicted bounding boxes to guide pixel-wise segmentation. Seo et al. [22] also present a framework called URVOS, which first predicts the initial mask based on the image and then utilizes the predicted mask of the previous frame for RVOS by memorizing the attention module.

Most recent RVOS works employ cross-attention to interact visual images with linguistic information. LBDT [6] uses language as an intermediate bridge to connect temporal and spatial information and leverages cross-modal attention operations to aggregate language-related motion and appearance. MMVT [40] calculates the optical flow between frames and fuses it as motion information with text features and visual features. However, these frame-based spatial granularity multimodal fusion methods have limitations and tend to lead to mismatches between visual and linguistic information. Therefore, a recent piece of work [32, 40] explores a novel multi-level representation learning method and introduces dynamic semantic alignment to adaptively fuse the two modal information.

Transformer [29] has been widely applied and achieved great success in many computer vision tasks, such as object detection [2, 43] and image segmentation [41, 5]. Since DETR [2] introduces a new query-based paradigm, the latest works [1, 33, 16] prefer to apply the DETR’s framework for RVOS task. Specifically, they utilize Transformer structures to interact visual images with linguistic data. Because Transformer has a remarkable ability to mine non-local correlation, they are able to attain SOTA performance in accuracy and efficiency. Despite the relative effectiveness of current RVOS techniques, they are primarily restricted to regular supervised learning settings, which would not be able to deal with unseen scenes with few shots.

### 3. Methods

#### 3.1. Overview

In the setting of FS-RVOS, we have training and testing datasets  $D_{train}$  and  $D_{test}$  with disjoint category sets  $C_{train}$  and  $C_{test}$ , i.e.,  $C_{train} \cap C_{test} = \emptyset$ . Similar to few-shot learning tasks [26], episode training is adopted in this work, where  $D_{train}$  and  $D_{test}$  consist of several episodes. Each episode contains a support set  $S$  and a query set  $Q$ , where the text-referred objects (target objects) from both sets belong to the same class. The support set has  $K$  image-mask pairs  $S = \{x_k, m_k\}_{k=1}^K$  and the corresponding referring expression with  $L$  words  $T_s = \{t_i\}_{i=1}^L$ , where  $m_k$  is the ground-truth mask of the video frame  $x_k$ . The query set  $Q = \{x_i^q\}_{i=1}^N$  is a selection of consecutive frames from a video and the corresponding natural language description with  $M$  words  $T_q = \{t_i\}_{i=1}^M$ ,  $N$  is the number of frames. With the setting above, FS-RVOS encourages models to segment objects with the unseen class in a query set based on a few samples in the support set.

As shown in Figure 2, our framework mainly consists of a Feature Extraction module, a Cross-modal Affinity module (CMA), and a Mask Generation module. With the support and query data as input, the framework predicts object masks for the query data, under the guidance of the corresponding language expressions. Specifically, the vision and text encoders extract features for visual and textual inputs, respectively. Then, CMA fuses visual and textual features hierarchically. The fused features are used to build the relationship between the support set and the query set. The output features are finally fed to the Mask Generation module to get the final segmentation results.

#### 3.2. Feature Extraction

We use a shared-weight visual encoder to extract multi-scale features from each frame in the support set and query set, resulting in visual feature sequences  $F_{vs} = \{f_{vs}\}_{vs=1}^K$  and  $F_{vq} = \{f_{vq}\}_{vq=1}^N$ . For linguistic information, we use a Transformer-based text encoder [19] to extract text features

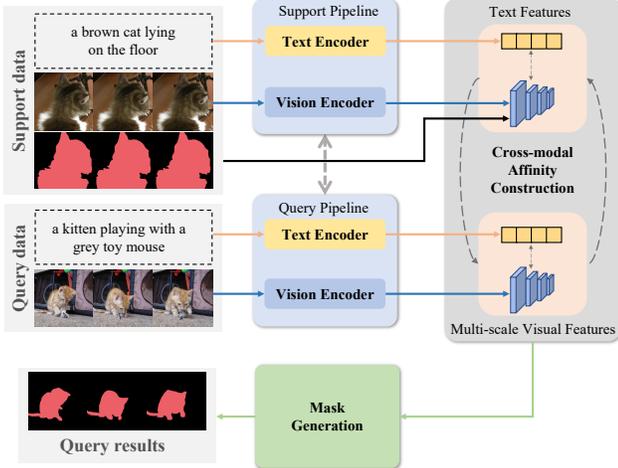


Figure 2: The overall pipeline of our framework. The feature encoder extracts visual and textual information from the support and query sets. Then, the Cross-modal Affinity module calculates the multimodal information affinity between the support set and the query set. Finally, the final segmentation result is obtained through Mask Generation.

$F_{ts} = \{f_{ts}\}_{ts=1}^L$  and  $F_{tq} = \{f_{tq}\}_{tq=1}^M$  from the natural language descriptions  $T_s$  and  $T_q$  that correspond to the input support data and the query data.

### 3.3. Cross-modal Affinity Construction

With a few samples in the support set, the goal of FS-RVOS is to efficiently leverage the given information and quickly adapt to relevant scenarios. Compared with conventional few-shot tasks, FS-RVOS not only builds the affinity between the support and query sets but also involves multimodal relationships between videos and referring expressions. Therefore, FS-RVOS is more challenging and requires specialized solutions for high-quality results.

We propose the Cross-Modal Affinity (CMA) module to achieve this, as shown in Figure 3. We first perform cross-attention fusion between the visual and textual features of support and query data, to obtain pixel-augmented multimodal features. Then, in order to aggregate beneficial information in support features, we build an affinity relationship between the support set and the query set.

Due to the diversity of referred objects and the drastic changes between video frames, it is a challenge to achieve the accurate location of the target only by visual information. Therefore, to obtain the accurate positioning of the segmentation target, we use language information as a supplement, which contains a specific description of the referred objects. To interact and align the visual features and text features, multi-head cross-attention (MCA) is proposed to fuse multimodal information, achieving two multi-scale

feature maps  $F'_{vs} = \{f'_s\}_{s=1}^K$  and  $F'_{vq} = \{f'_q\}_{q=1}^N$ :

$$f'_{vs} = \text{MCA}(f_{vs}, f_{ts}), f'_{vq} = \text{MCA}(f_{vq}, f_{tq}), \quad (1)$$

where  $f_{vs}, f_{vq}$  represent the visual features of support and query,  $f_{ts}, f_{tq}$  are their corresponding textual features. Here we calculate an affinity between textual features and visual features to filter out irrelevant visual information. Compared with concatenation, MCA suits our framework better since it can leverage the similarities between multimodal features for information complementation.

The affinity between the support set and query set indicates the multimodal feature correlation among them, providing valuable clues for the segmentation of the query data. Although the objects to segment in the support and query sets belong to the same category, their visual properties usually have significant differences, such as appearance, pose, and scene. This means only a tiny part of the information in the support data is conducive to segmenting the query data, while other information will cause bad results. Therefore, we urgently need to solve the problem of computing the correct affinity relationship between multimodal information in the support and query sets.

To achieve this, we propose a self-affinity block to encode the query features and a cross-affinity block to enable the query features to focus on beneficial pixels in the support features. In particular, given the input query features, we utilize a convolution operation to map them as query  $q_q$ , key  $k_q$ , and value  $v_q$ . We perform the same operation for the support features to map them to key  $k_s$  and value  $v_s$ . The input of the self-affinity block does not include the support features and mainly aggregates the context information of the query features for better segmentation.

First, we calculate the affinity map  $A^Q = \frac{q_q \cdot (k_q)^T}{\sqrt{d_{head}}}$ , where  $d_{head}$  is the hidden dimension of the input sequences, and we assume all sequences have the same dimension 256 by default. Therefore, the query features after the self-affinity block are represented as:

$$q_s = \text{Softmax}(A^Q)v_q. \quad (2)$$

We then feed the obtained query features to the cross-affinity block. The purpose of the cross-affinity block is to construct the cross-affinity relationship between the support features and the query features and aggregate the useful information. Our cross-affinity block can be formulated as:

$$query_{feat} = \text{Softmax}\left(\frac{q_s \cdot (k_s)^T}{\sqrt{d_{head}}}\right)v_s, \quad (3)$$

where  $q_s$  is the output of the self-affinity block. Through these two modules, query features enhance features by modeling contextual information and computing the correlation between support and query features, effectively avoiding the attention bias caused by irrelevant features.

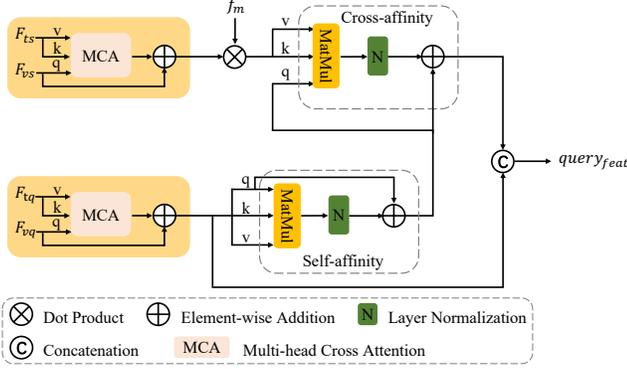


Figure 3: The architecture of the Cross-modal Affinity (CMA) module. We use multi-head cross-attention to fuse visual and text features to get more robust features. Self-affinity for modeling contextual information on query features and cross-affinity for aggregating beneficial information from support features.

### 3.4. Mask Generation

The goal of the Mask Generation module is to find the most relevant objects and decode the features step by step. To achieve this goal, the structure of Deformable-DETR [43] and feature pyramid [17] has been used in our work. We add the corresponding positional encoding to the feature sequence that aggregates beneficial information, which is then sent to the Transformer encoder. In the decoder part of Transformer, we introduce  $N$  learnable anchor boxes as queries to represent the instances of each frame. These queries are replicated as decoder input for all frames and finally converted into instance embeddings by the decoder, resulting in  $N_q = T \times N$  predictions.

In the feature pyramid network, in order to gradually fuse multimodal features from different layers, the output of the Transformer encoder and the features from the vision encoder are stacked to form hierarchical features. We use  $f_v^l$  to represent visual features at each level. First, we down-sample the multi-scale visual features, and the time dimension remains unchanged. Then interact with visual and linguistic features in a cross-attention manner to enhance object pixel features, thereby facilitating mask prediction. The fused features are upsampled to restore the previous shape:

$$Cross(f_v^l, f_{tq}) = \text{Softmax}\left(\frac{f_v^l \cdot (f_{tq})^T}{\sqrt{d_{head}}}\right) f_{tq}, \quad (4)$$

where  $f_{tq}$  represents the text features corresponding to the query set. Finally, we pass the features of the last layer through a  $3 \times 3$  convolutional layer to get the final feature maps  $F_{seg} = \{f_{seg}^t\}_{t=1}^T$ , where  $f_{seg}^t \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C}$ .

Furthermore, we construct a set of auxiliary heads to obtain the final object mask across frames. The class head is used to calculate the confidence score  $s_t$ . The score indi-

cates whether the instance corresponds to the referred object and whether the object is visible in the current frame. The kernel head is implemented by three consecutive linear layers, which generate the parameters  $W = \{w_t\}_{t=1}^{N_q}$  of the  $N_q$  dynamic kernels. We use them as convolution filters on the feature maps, generating a series of segmentation masks.

### 3.5. Loss Function

In Mask Generation, we use  $N$  learnable anchor boxes as queries and generate a set of  $N_q = T \times N$  prediction sequences. We denote the predicted sequence as  $\hat{y} = \{\hat{y}_i\}_{i=1}^{N_q}$  and the prediction for the  $i^{\text{th}}$  instance is expressed as:

$$\hat{y}_i = \{\hat{s}_i^t, \hat{m}_i^t\}_{t=1}^T, \quad (5)$$

where  $\hat{s}_i^t \in \mathbb{R}^1$  is a fractional score indicating whether the instance corresponds to the referred object.  $\hat{m}_i^t \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4}}$  is the predicted binary segmentation mask. Similar to the previous methods [1, 33], we use dynamic convolution to generate the object mask. We get the final feature maps  $f_{seg}^t$  through the feature pyramid network, and the mask prediction can be calculated by  $\hat{m}^t = \{w_t * f_{seg}^t\}$ .

The sequence of ground-truth instances is denoted as  $y = \{s^t, m^t\}_{t=1}^T$ , and  $s^t$  is a one-hot value and equals 1 when the ground-truth instance is visible in the frame, and 0 otherwise. Our loss function is defined as:

$$\mathcal{L}(y, \hat{y}_i) = \lambda_{cls} \mathcal{L}_{cls}(y, \hat{y}_i) + \lambda_{kernel} \mathcal{L}_{kernel}(y, \hat{y}_i), \quad (6)$$

where  $\lambda_{cls}$ ,  $\lambda_{kernel}$  are hyperparameters to balance the loss. We use focal loss [18] (denoted as  $\mathcal{L}_{cls}$ ) to supervise the prediction of instance sequence reference results.  $\mathcal{L}_{kernel}$  is a combination of DICE loss [20] and the binary mask focal loss.

### 3.6. Inference

In the inference phase, we treat all input video frames as a whole and predict the mask trajectory for the entire video using only one forward pass. Given an input video and the corresponding linguistic expression, our model generates a sequence of  $N$  instances. For each frame, we select the instance sequence with the highest confidence score as the final prediction, and its index can be expressed as:

$$\bar{s} = \arg \max_{i \in \{1, 2, \dots, N\}} s_i. \quad (7)$$

The final prediction  $m = \{s_t\}_{t=1}^T$  for each frame is obtained from the mask candidate set  $\hat{m}^t$  with index  $\bar{s}$ .

## 4. Benchmark

Since the existing RVOS datasets [22, 14, 34, 11] only target specific scenarios, the model cannot handle the diverse scenarios in the real world. Likewise, these datasets



A man in a brown jacket is standing on the right



A man in a gray coat standing on the right is saying something



A brown dog is barking

Figure 4: Annotation examples of the Mini-Ref-SAIL-VOS dataset.

are not suitable for the few-shot RVOS problem. This is because the train/test/validation sets of these datasets have class repetition and cannot be used to evaluate the generality of unseen classes.

**Mini-Ref-YouTube-VOS.** To match the FS-RVOS setting, we build up a new dataset called Mini-Ref-YouTube-VOS based on the Ref-YouTube-VOS dataset [22]. The data that can be directly obtained from the Ref-YouTube-VOS dataset contains 3,471 videos, 12,913 referring expressions, and annotated instances covering more than 60 categories. However, some videos in this dataset consist of multiple category instances. When preparing data for the few-shot setting, we cleaned up the dataset, i.e., removing such videos and keeping only those containing only one category instance, a total of 2387 videos were obtained.

The video data in the dataset should be class-balanced, and the number of samples in each class should not vary too much to avoid overfitting any class. Therefore, we deal with the categories whose number of videos does not meet the requirements to ensure the class balance of the dataset. After the above screening, 1,668 videos were obtained, including 48 classes. To better show the model results, we adopt the cross-validation method to divide the dataset into four folds on average. Each fold contains 36 training and 12 test classes with disjoint categories.

**Mini-Ref-SAIL-VOS.** Most Mini-Ref-YouTube-VOS data involve natural scenes of a relatively homogeneous type and therefore do not represent the diversity of data in the real world. To better demonstrate the generalization of our model, we collect videos from SAIL-VOS [13] to

construct a new dataset Mini-Ref-SAIL-VOS. The SAIL-VOS dataset is a synthetic dataset collected from a video game GTA-V, aiming to foster semantic amodal segmentation research. In SAIL-VOS, each frame is accompanied by densely annotated, pixel-wise and amodal segmentation masks with semantic labels. Since the data is collected from the game, phenomena such as shot transition, segmentation target cross-frame, and target occlusion are inevitable, which brings challenges to the segmentation task.

We reorganize the SAIL-VOS dataset to pick out samples suitable for the FS-RVOS setting. First, for videos with few frames for segmentation targets, we directly discard them. For the case where the segmentation target appears across frames, we manually delete the frames where the target does not appear in the middle to maintain the temporal continuity of the segmentation target. For phenomena such as object occlusion, we characterize it as a challenge and do not deliberately delete video frames with object occlusion. Following the above settings, we collected a dataset with 68 videos and 3 semantic categories.

It is worth noting that although there are accurate mask annotations in the SAIL-VOS dataset, natural language description corresponding to the segmentation target is not available. Thus, to adapt it to FS-RVOS, we employed expert annotators to provide referring expressions after data collection. Given a pair of videos for each annotator, the video frames are superimposed with corresponding masks to indicate the objects to be segmented. The annotators were then asked to provide a distinguishing statement with a word limit of 20 words. To ensure the quality of natural language annotations, all annotations are verified and cleaned

	Method	Fold-1	Fold-2	Fold-3	Fold-4	Mean
$\mathcal{J}$	DANet [3]	47	33.5	38.5	44	40.8
	Ours	<b>59.5</b>	<b>45.3</b>	<b>50.4</b>	<b>57.3</b>	<b>53.1</b>
$\mathcal{F}$	DANet [3]	49.3	38.2	41.4	45.8	43.7
	Ours	<b>60.8</b>	<b>48.9</b>	<b>51.3</b>	<b>58.1</b>	<b>54.8</b>

Table 1: Quantitative results on Mini-Ref-YouTube-VOS. We added visual language fusion modules to DANet.

	Method	Fold-1	Fold-2	Fold-3	Fold-4	Mean
$\mathcal{J}$	DANet [3]	54.3	47.6	30.9	30.6	40.9
	Ours	<b>80.9</b>	<b>80.8</b>	<b>80.1</b>	<b>68.9</b>	<b>77.7</b>
$\mathcal{F}$	DANet [3]	54.1	48.4	35.2	36	43.4
	Ours	<b>77.1</b>	<b>77</b>	<b>77.3</b>	<b>67.8</b>	<b>74.8</b>

Table 2: Quantitative results on Mini-Ref-SAIL-VOS.

up after the initial annotation. The target will not be used if the natural language description cannot clearly describe the target. As shown in Figure 4, we show the selected videos along with referring expressions.

## 5. Experiments

### 5.1. Implementation Details

We adopt ResNet-50 [12] and RoBERTa-Base [19] as our vision and text encoders, respectively. During the training stage, the parameters of both encoders are frozen. In our experiments, we adopt a 5-shot setting. Specifically, we extract 5 consecutive frames and the corresponding referring expressions from a certain video of a class as a support set. The query set is composed of consecutive frames and corresponding referring expressions extracted from other videos belonging to the same class. We test each fold 5 times and report the average confidence of the results. Our model utilizes AdamW for optimization. The weight decay is  $5 \times 10^{-4}$  and the initial learning rate is  $1 \times 10^{-4}$ . To balance GPU memory efficiency, we downsample all video frames, with the shortest video frame size being 360 and the longest 640. The parameters of the loss function are set as  $\lambda_{cls} = 2$ ,  $\lambda_{kernel} = 5$ . All methods for conducting experiments will be pre-processed and fine-tuned in the same way, i.e., pre-trained on Ref-COCO [36] dataset.

Following the settings of previous RVOS works, we use the region similarity ( $\mathcal{J}$ ) and the contour accuracy ( $\mathcal{F}$ ) to measure the model performance.

### 5.2. Results

As a novel problem, no relevant works can be directly used for comparisons. Therefore, we choose DANet [3] as the baseline given we both focus on few-shot video segmentation. For a fair comparison, we add a visual-language fusion module to the Few-Shot VOS model.

Method	$\mathcal{J}$	$\mathcal{F}$	$\mathcal{J}\&\mathcal{F}$
LBDT [6]	27.5 / <u>42.4</u>	36.2 / <u>37.3</u>	31.6 / <u>39.6</u>
ReferFormer [33]	65.1 / <u>74.1</u>	62.8 / <u>64.9</u>	64.0 / <u>69.5</u>
MTTR [1]	66.5 / <u>69.7</u>	64.9 / <u>68.1</u>	65.7 / <u>68.9</u>
Ours	<b>80.9</b>	<b>77.1</b>	<b>79</b>

Table 3: Comparison with state-of-the-art methods from RVOS on the Mini-Ref-SAIL-VOS dataset to measure the model’s generalization. Underlined scores are achieved after fine-tuning.

Self-affinity	Cross-affinity	$\mathcal{J}\&\mathcal{F}$
-	-	57.9
✓	-	59.2
✓	✓	60.2

Table 4: Ablation studies that validate the effectiveness of each component in our CMA. The first result is obtained with our baseline.

**Mini-Ref-YouTube-VOS.** We present the experimental results of our model on the Mini-Ref-YouTube-VOS dataset in Table 1. It can be observed that our method significantly outperforms DANet. Compared with previous methods, our method achieves a substantial increase in average performance, with an average improvement of more than 10%. The excellent performance demonstrates the superiority and robustness of our proposed method.

**Mini-Ref-SAIL-VOS.** To evaluate the generalization of our model, we make further experiments and comparisons on the Mini-Ref-SAIL-VOS dataset. We do not perform new training on the Mini-Ref-SAIL-VOS dataset but directly test with the model trained on Mini-Ref-YouTube-VOS. Note that the videos in the Mini-Ref-SAIL-VOS dataset are from Game scenes, and they hold a noticeable domain different from the data in the Mini-Ref-YouTube-VOS dataset. In addition, the objects in some videos are occluded and these phenomena make this dataset somewhat challenging. We show the experimental results in Table 2. According to the Table, it is clear that our method achieves a significant improvement over the baseline.

Moreover, we also make further comparisons with some state-of-the-art RVOS methods [6, 33, 1]. Here, we only show the results under Fold-1. To measure the RVOS model’s generalization, we first directly test the trained RVOS models on the Mini-Ref-SAIL-VOS dataset. Since our few-shot learning task needs a few samples as support data, for a fair comparison, these RVOS models will be fine-tuned with a few samples. The corresponding results are shown in Table 3. From the table, it can be seen that although the fine-tuning of the existing models can lead to performance gains, there is still a big gap compared to our

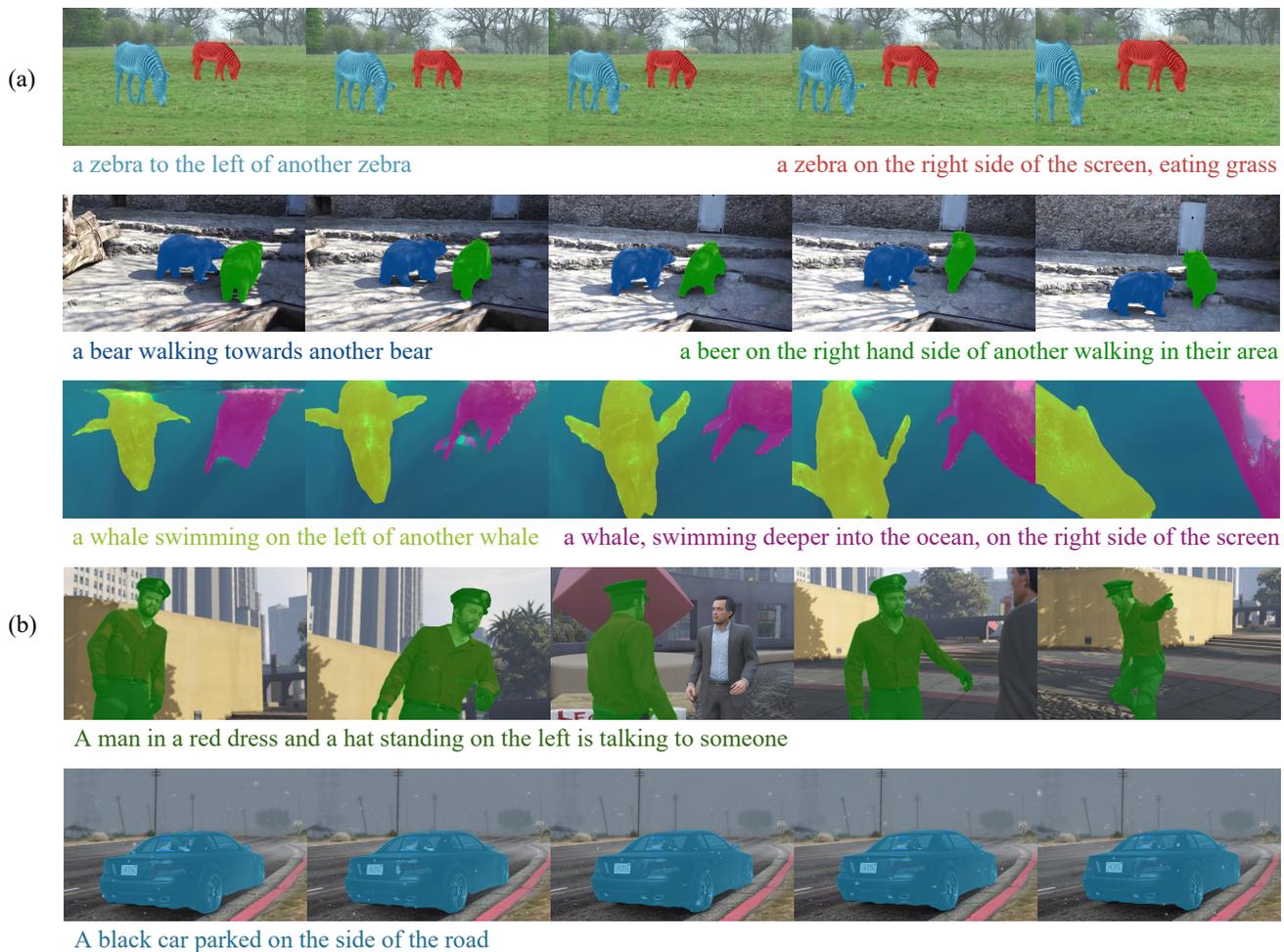


Figure 5: Qualitative results on (a) Mini-Ref-YouTube-VOS and (b) Mini-Ref-SAIL-VOS.

proposed approach. This is because our model effectively constructs the multimodal relationship between the support set and the query set so that it can quickly adapt to new scenarios by only using a few samples.

### 5.3. Ablation Study

In this section, we perform an ablation study on the Mini-Ref-YouTube-VOS dataset to evaluate the design and the robustness of the model. Unless otherwise stated, we only show results under Fold-1. Denoting  $\mathcal{J}$  and  $\mathcal{F}$  as the average of  $\mathcal{J}$  and  $\mathcal{F}$ , we can use the indicator to show the performance of the model.

**Cross-modal Affinity.** We perform ablation studies on the components of the CMA in Table 4. The results of the baseline are shown in the first line. As mentioned before, the baseline refers to directly concatenating the support features and query features into the Mask Generation to obtain the segmentation mask.

First, we only utilize the self-affinity block to establish contextual information between pixels to enhance query

features. At this time, support features are concatenated with the enhanced query features and then fed to the Mask Generation module. It can be seen that good results are achieved, indicating that the Transformer does work for modeling features and extracting contextual information. By adding the proposed cross-affinity block, the performance can be further improved by 1%. Such comparisons show that the cross-affinity block properly constructs the multimodal relationship between the support set and the query set, effectively avoiding the bias of the query features.

### 5.4. Qualitative Results

The qualitative results of our model are presented in Figure 5. From the figure, It can be seen that the proposed model can segment the referred objects accurately in a variety of challenging situations. Furthermore, we also show the qualitative results on the Mini-Ref-SAIL-VOS dataset. In general, our model always achieves high-quality results even in the face of samples from different scenes.

## 6. Conclusion

In this work, we propose CMA to learn multimodal affinity in a few samples to segment diverse data. Further, we generalize it as a few-shot RVOS problem. We validate our model on the newly constructed datasets - Mini-Ref-YouTube-VOS and Mini-Ref-SAIL-VOS and obtain state-of-the-art performance. We hope this work can cast a light on future FS-RVOS research.

**Acknowledgements.** This work is supported in part by the National Key R&D Program of China (Grant No. 2022YFF1202903), the National Natural Science Foundation of China (Grant No. 61971004, 62122035), the Natural Science Foundation of Anhui Province, China (Grant No. 2008085MF190), and the Equipment Advanced Research Sharing Technology Project, China (Grant No. 80912020104).

## References

- [1] Adam Botach, Evgenii Zheltonozhskii, and Chaim Baskin. End-to-end referring video object segmentation with multimodal transformers. *arXiv preprint arXiv:2111.14821*, 2021. 2, 3, 5, 7
- [2] 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
- [3] Haoxin Chen, Hanjie Wu, Nanxuan Zhao, Sucheng Ren, and Shengfeng He. Delving deep into many-to-many attention for few-shot video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14040–14049, 2021. 3, 7
- [4] Junjie Chen, Li Niu, Siyuan Zhou, Jianlou Si, Chen Qian, and Liqing Zhang. Weak-shot semantic segmentation via dual similarity transfer. In *NeurIPS*, 2022. 3
- [5] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*, 34, 2021. 3
- [6] Zihan Ding, Tianrui Hui, Junshi Huang, Xiaoming Wei, Jizhong Han, and Si Liu. Language-bridged spatial-temporal interaction for referring video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4964–4973, 2022. 3, 7
- [7] Nanqing Dong and Eric P Xing. Few-shot semantic segmentation with prototype learning. In *BMVC*, volume 3, 2018. 2
- [8] Qi Fan, Wenjie Pei, Yu-Wing Tai, and Chi-Keung Tang. Self-support few-shot semantic segmentation. 2022. 2
- [9] Tsu-Jui Fu, Xin Eric Wang, Scott T Grafton, Miguel P Eckstein, and William Yang Wang. M3l: Language-based video editing via multi-modal multi-level transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10513–10522, 2022. 1
- [10] Mingqi Gao, Feng Zheng, James JQ Yu, Caifeng Shan, Guiguang Ding, and Jungong Han. Deep learning for video object segmentation: a review. *Artificial Intelligence Review*, 56(1):457–531, 2023. 1
- [11] Kirill Gavriluk, Amir Ghodrati, Zhenyang Li, and Cees GM Snoek. Actor and action video segmentation from a sentence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5958–5966, 2018. 3, 5
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 7
- [13] Yuan-Ting Hu, Hong-Shuo Chen, Kexin Hui, Jia-Bin Huang, and Alexander G Schwing. Sail-vos: Semantic amodal instance level video object segmentation-a synthetic dataset and baselines. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3105–3115, 2019. 2, 6
- [14] Anna Khoreva, Anna Rohrbach, and Bernt Schiele. Video object segmentation with language referring expressions. In *Asian Conference on Computer Vision*, pages 123–141. Springer, 2018. 3, 5
- [15] Chunbo Lang, Gong Cheng, Binfei Tu, and Junwei Han. Learning what not to segment: A new perspective on few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8057–8067, June 2022. 3
- [16] Xiang Li, Jinglu Wang, Xiaohao Xu, Xiao Li, Yan Lu, and Bhiksha Raj. R<sup>2</sup>vos: Robust referring video object segmentation via relational multimodal cycle consistency. *arXiv preprint arXiv:2207.01203*, 2022. 3
- [17] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 5
- [18] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 5
- [19] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 3, 7
- [20] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. IEEE, 2016. 5
- [21] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 622–631, 2019. 2
- [22] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object segmentation network with a large-scale benchmark. In *European Conference on Computer Vision*, pages 208–223. Springer, 2020. 2, 3, 5, 6

- [23] Amirreza Shaban, Shray Bansal, Zhen Liu, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. In *BMVC*, 2017. 2
- [24] Mennatullah Siam, Konstantinos G Derpanis, and Richard P Wildes. Temporal transductive inference for few-shot video object segmentation. *arXiv preprint arXiv:2203.14308*, 2022. 3
- [25] Mennatullah Siam, Naren Doraiswamy, Boris N Oreshkin, Hengshuai Yao, and Martin Jagersand. Weakly supervised few-shot object segmentation using co-attention with visual and semantic embeddings. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 860–867, 2021. 3
- [26] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. 3
- [27] Zhuotao Tian, Xin Lai, Li Jiang, Shu Liu, Michelle Shu, Hengshuang Zhao, and Jiaya Jia. Generalized few-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11563–11572, June 2022. 2
- [28] Zhuotao Tian, Hengshuang Zhao, Michelle Shu, Zhicheng Yang, Ruiyu Li, and Jiaya Jia. Prior guided feature enrichment network for few-shot segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 2
- [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 3
- [30] Hao Wang, Cheng Deng, Fan Ma, and Yi Yang. Context modulated dynamic networks for actor and action video segmentation with language queries. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12152–12159, 2020. 3
- [31] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9197–9206, 2019. 2
- [32] Dongming Wu, Xingping Dong, Ling Shao, and Jianbing Shen. Multi-level representation learning with semantic alignment for referring video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4996–5005, 2022. 3
- [33] Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4974–4984, 2022. 2, 3, 5, 7
- [34] Chenliang Xu, Shao-Hang Hsieh, Caiming Xiong, and Jason J Corso. Can humans fly? action understanding with multiple classes of actors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2264–2273, 2015. 5
- [35] Boyu Yang, Chang Liu, Bohao Li, Jianbin Jiao, and Qixiang Ye. Prototype mixture models for few-shot semantic segmentation. In *European Conference on Computer Vision*, pages 763–778. Springer, 2020. 2
- [36] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *European Conference on Computer Vision*, pages 69–85. Springer, 2016. 7
- [37] Chi Zhang, Guosheng Lin, Fayao Liu, Jiushuang Guo, Qingyao Wu, and Rui Yao. Pyramid graph networks with connection attentions for region-based one-shot semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9587–9595, 2019. 2
- [38] Gengwei Zhang, Guoliang Kang, Yi Yang, and Yunchao Wei. Few-shot segmentation via cycle-consistent transformer. *Advances in Neural Information Processing Systems*, 34:21984–21996, 2021. 2
- [39] Shan Zhang, Tianyi Wu, Sitong Wu, and Guodong Guo. Catrans: Context and affinity transformer for few-shot segmentation. *arXiv preprint arXiv:2204.12817*, 2022. 2
- [40] Wangbo Zhao, Kai Wang, Xiangxiang Chu, Fuzhao Xue, Xinchao Wang, and Yang You. Modeling motion with multi-modal features for text-based video segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11737–11746, 2022. 3
- [41] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. 3
- [42] Siyuan Zhou, Li Niu, Jianlou Si, Chen Qian, and Liqing Zhang. Weak-shot semantic segmentation by transferring semantic affinity and boundary. *arXiv preprint arXiv:2110.01519*, 2021. 3
- [43] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3, 5