# Language as Queries for Referring Video Object Segmentation

Jiannan Wu[1],    Yi Jiang[2]    Peize Sun[1],    Zehuan Yuan[2],    Ping Luo[1]

[1]The University of Hong Kong    [2]ByteDance

## Abstract

*Referring video object segmentation (R-VOS) is an emerging cross-modal task that aims to segment the target object referred by a language expression in all video frames. In this work, we propose a simple and unified framework built upon Transformer, termed ReferFormer. It views the language as queries and directly attends to the most relevant regions in the video frames. Concretely, we introduce a small set of object queries conditioned on the language as the input to the Transformer. In this manner, all the queries are obligated to find the referred objects only. They are eventually transformed into dynamic kernels which capture the crucial object-level information, and play the role of convolution filters to generate the segmentation masks from feature maps. The object tracking is achieved naturally by linking the corresponding queries across frames. This mechanism greatly simplifies the pipeline and the end-to-end framework is significantly different from the previous methods. Extensive experiments on Ref-Youtube-VOS, Ref-DAVIS17, A2D-Sentences and JHMDB-Sentences show the effectiveness of ReferFormer. On Ref-Youtube-VOS, ReferFormer achieves 55.6 $\mathcal{J}\&\mathcal{F}$ with a ResNet-50 backbone without bells and whistles, which exceeds the previous state-of-the-art performance by 8.4 points. In addition, with the strong Swin-Large backbone, ReferFormer achieves the best $\mathcal{J}\&\mathcal{F}$ of 64.2 among all existing methods. Moreover, we show the impressive results of 55.0 mAP and 43.7 mAP on A2D-Sentences and JHMDB-Sentences respectively, which significantly outperforms the previous methods by a large margin. Code is publicly available at https://github.com/wjn922/ReferFormer.*

## 1. Introduction

Referring video object segmentation (R-VOS) aims to segment the target object in a video given a natural language description. This emerging topic has raised great attention in the research community and is expected to benefit many applications in a friendly and interactive way, *e.g.*, video editing and video surveillance. R-VOS is more challenging than the traditional semi-supervised video object segmenta-
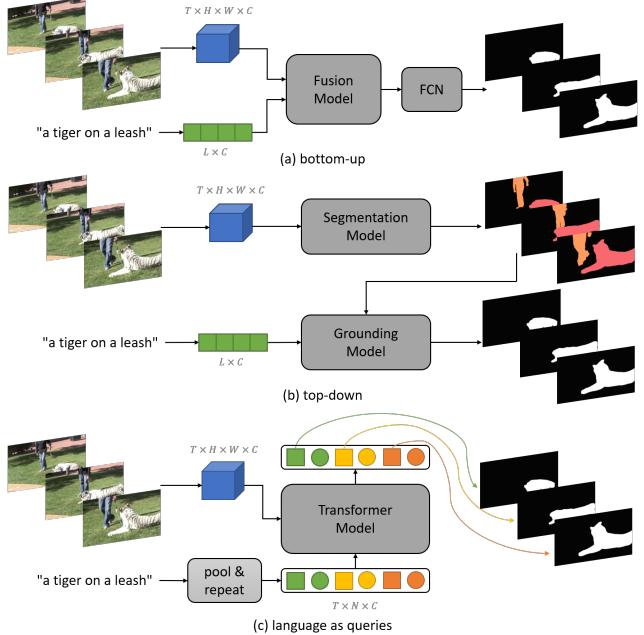


Figure 1. Comparison of current referring video object segmentation (R-VOS) pipelines. **(a)** Bottom-up. **(b)** Top-down. **(c)** Ours.

tion [39, 53], because it does not only lack the ground-truth mask annotation in the first frame, but also require the comprehensive understanding of the cross-modal sources, *i.e.*, vision and language. Therefore, the model should have a strong ability to infer which object is referred and to perform accurate segmentation.

To accomplish this task, the existing methods can be mainly categorized into two groups: (1) *Bottom-up* methods. These methods incorporate the vision and language features in a early-fusion manner, and then adopt a FCN [32] as decoder to generate object masks, as shown in Figure 1(a). (2) *Top-down* methods. These methods tackle the problem in a top-down perspective and follow a two-stage pipeline. As illustrated in Figure 1(b), they first employ an instance segmentation model to find all the objects in each frame, and then associate them in the entire video to form the tracklet candidates. Afterwards, they use the expression as the grounding criterion to select the best-matched one.

Although these two streams of methods have demon-

strated their effectiveness with promising results, they still have some intrinsic limitations. First, for the *bottom-up* methods, they fail to capture the crucial instance-level information and do not consider the object association across multiple frames. Therefore, this type of methods can not provide explicit knowledge for cross-modal reasoning and would encounter the discrepancy of predicted object due to scene changes. Second, although *top-down* methods have greatly boost the performance over the *bottom-up* methods, they suffer from heavy workload because of the complex, multi-stage pipeline. For example, the recent method proposed by Liang *et al*. [23] comprises of three parts: HTC [5], CFBI [56] and a tracklet-language grounding model. All these networks need to be pretrained on the ImageNet [21], COCO [27] or RefCOCO [59] and further finetuned on R-VOS datasets, respectively. Furthermore, the separate optimization on several sub-problems would lead to sub-optimal solution.

These limitations of current methods motivate us to design a simple and unified framework that solves the R-VOS task elegantly. The recent success of Transformer [46] in object detection [4, 62] and video instance segmentation [17, 50, 51] demonstrates a promising solution. However, it is non-trivial to apply such models to the R-VOS task. These models [4,62] use a fixed number (*e.g*., 100) of learnable queries to detect all the objects in an image. Under this circumstance, it would be confused for the model to distinguish which object is referred due to the randomness of the expression. Here raises a natural question: *"Is it possible for a unified model to know where to look using queries?"*

This work answers the question by proposing the notion of *language as queries*, as shown in Figure 1(c). We put the linguistic restriction on all object queries and use these *conditional queries* as input for the model. In this manner, the expression will make the queries focus on the referred object only, and thus greatly reducing the query number (*e.g*., 5 in our experiments). The next challenge lies in how to decode the object mask from query representations. As the queries contain rich instance characteristics, we view them as instance-aware dynamic kernels to filter out the segmentation masks from feature maps. Moreover, to make the feature maps more discriminative, we design a novel cross-modal feature pyramid network (CM-FPN) where the visual and linguistic features interact in multiple levels for fine-grained cross-modal fusion.

The unified framework can not only produce the segmentation masks for referred objects, but also the classification results and detection boxes. Moreover, the *conditional queries* are linked via instance matching strategy across frames so that the object tracking is achieved naturally without post-process. As shown in Figure 5, our unified framework is able to detect, segment and track the referred object simultaneously. We hope this framework could serve as a strong baseline for R-VOS task.

The main contributions of this work are as follows.

- We propose a simple and unified framework for referring video object segmentation, termed RefeFormer. Given a video clip and the corresponding language expression, our framework directly detects, segments and tracks the referred object in all frames in an end-to-end manner.

- We present the notion of *language as queries*. We introduce a small set of object queries which conditioned on the text expression to attend the referred object only. These *conditional queries* are shared across different frames in the initial state and they are transformed into dynamic kernels to filter out the segmentation masks from feature maps. This mechanism provides a new perspective for the R-VOS task.

- We design the cross-modal feature pyramid network (CM-FPN) for multi-scale vision-language fusion, which improves the discriminativeness of mask features for accurate segmentation.

- Extensive experiments on Ref-Youtube-VOS, Ref-DAVIS17, A2D-Sentences and JHMDB-Sentences show that RefeFormer outperforms the previous methods on these four benchmarks by a large margin. *E.g*., on Ref-Youtube-VOS, RefeFormer with a ResNet-50 backbone achieves 55.6 $\mathcal{J}\&\mathcal{F}$ without bells and whistles, showing the significant 8.4 points gain over the previous state-of-the-art methods. And using the strong Video-Swin-Base visual backbone, RefeFormer achieves the impressive results of 64.9 $\mathcal{J}\&\mathcal{F}$.

## 2. Related Work

**Semi-supervised Video Object Segmentation.** The traditional semi-supervised video object segmentation (Semi-VOS) aims to propagates the ground-truth object masks given in the first frame to the entire video. Most recent works [8, 38, 47, 56] lie in the group of matching-based methods, which perform feature matching to track the target objects. STM [38] leverages a memory to store the past object features and utilize the attention matching mechanism on the memory to guide the prediction of current frame. CFBI [56] not only considers the embedding learning of foreground objects but also the background, resulting in a more robust framework.

**Referring Video Object Segmentation.** Referring video object segmentation (R-VOS) provides the language description instead of mask annotation as the object reference, thus it would be a more challenging task. The current methods for R-VOS mainly follow the two pipelines: (1) *Bottom-up* methods. An intuitive thinking is directly applying the image-level methods [9, 16, 34, 60] on the video frames independently, *e.g*., RefVOS [2]. The obvious drawback of such methods is that they fail to utilize
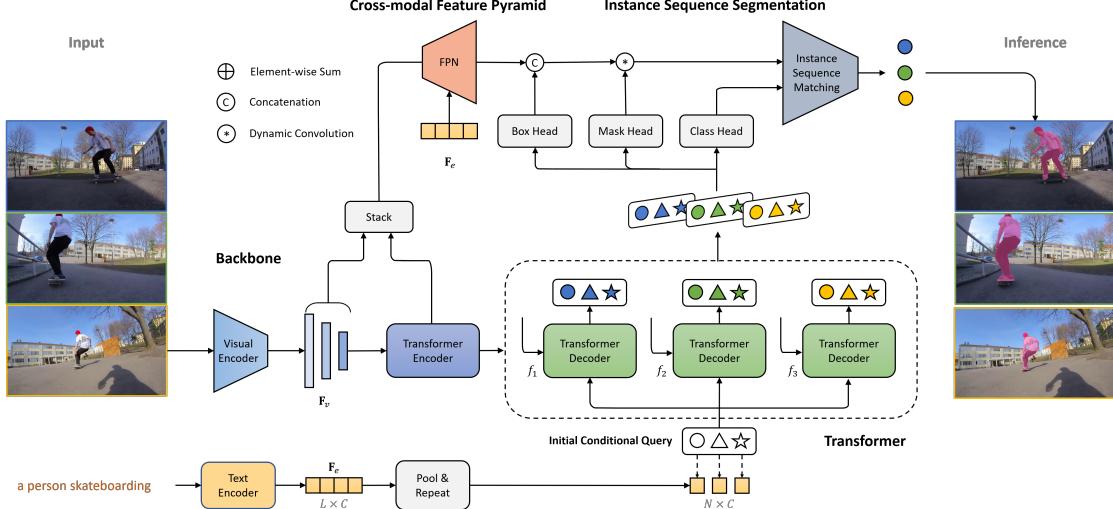
Figure 2. The overall pipeline of ReferFormer. It mainly consists of four parts: Backbone, Transformer, Cross-modal Feature Pyramid and the Segmentation part. The model takes a video clip with the corresponding language expression as input and output the segmentation mask of the referred object in each frame. For the Transformer decoder input, the object queries are conditioned on the language expression to find the referred object. The same colors represent the queries in the same frame and the same shapes represent the queries refer to the same instance. The order of queries inner frame keep consistent for different frames. Best viewed in color.

the valuable temporal information across frames, resulting in inconsistent object prediction due to the scene or appearance variations. To address this issue, URVOS [42] casts the task as a joint problem of referring object segmentation in an image and mask propagation in a video. They propose a unified referring VOS framework that employs a memory attention module to leverage the information of mask predictions in previous frames. (2) *Top-down* methods. The typical top-down method [23] first constructs an exhaustive set of object tracklets by propagating the object masks detected from several key frames to the whole video. Then, a language grounding model is built to select the best object tracklet from the candidate set. Although the method has made breakthrough performance improvement over the previous methods, the complex, multi-stage pipeline is computational-expensive and impractical.

In contrast to these two pipelines, we propose a *query-based* method that achieves the strongest performance with a simple and unified framework. The very recent work MTTR [3] also relies on the query-based mechanism. Nevertheless, they need the exhaustive segmentation annotations of all objects and supervised the un-referred instances during training process, which increases the workload of laborious annotation and makes the framework limited in practical applications.

**Transformer** Transformer [46] was first introduced for sequence-to-sequence translation in natural language processing (NLP) community and has achieved marvelous success in most computer vision tasks [11, 19, 30] such as object detection [4, 62], tracking [6, 36, 43, 54] and segmentation [7, 17, 61]. DETR [4] introduces the new *query-based*

paradigm [44, 62] for object detection, which employs a set of object queries as candidates and inputs them to the Transformer decoder. Beyond image field, VisTR [50] extends the framework for video instance segmentation (VIS) [55] task and solves the problem in a direct end-to-end parallel sequence decoding manner. SeqFormer [51] decouples the content query and box query to aggregates temporal information from each frame and achieves the state-of-the-art performance on VIS task. Inspired by these works, our work also relies on the *query-based* mechanism of Transformer but considers an additional modality, *i.e.*, language, as the object reference. Thus, we propose the notion of *language as queries* and build the simple and unified framework that detects, segments and tracks the referred object simultaneously.

## 3. Approach

Given a video clip $\mathcal{I} = \{I_t\}_{t=1}^{T}$ with $T$ frames and a referring expression $\mathcal{E} = \{e_l\}_{l=1}^{L}$ with $L$ words, we aim to produce $T$-frame binary segmentation masks of referred object $\mathcal{S} = \{s_t\}_{t=1}^{T}, s_t \in \mathbb{R}^{H \times W}$ in an end-to-end manner. To this end, we propose a simple and unified framework named ReferFormer, as shown in Figure 2. It mainly consists of four key components: Backbone, Transformer, Cross-modal Feature Pyramid network (CM-FPN) and the Instance Sequence Segmentation process. A small set of object queries conditioned on the language is introduced to find the referred object. During inference, we directly output the mask predictions by selecting the queries with the highest average score as the final results.

## 3.1. Backbone

**Visual Encoder.** We start by adopting a visual backbone to extract the multi-scale feature maps for each frame in the video clip independently, resulting in the visual feature sequence $\mathcal{F}_v = \{f_t\}_{t=1}^T$. It is noteworthy that both the 2D spatial encoder (*e.g.*, ResNet [14]) and 3D spatio-temporal encoder (*e.g.*, Video Swin Transformer [31]) could play the role of visual backbone.

**Linguistic Encoder.** Given the language description with $L$ words, we use off-the-shelf linguistic embedding model, RoBERTa [29], to extract the text feature $\mathcal{F}_e = \{f_i\}_{i=1}^L$. And we also obtain the sentence-level feature $f_e^s \in \mathbb{R}^C$ by pooling the features of each word. They are both necessary and essential in our model, because the sentence feature guides the learnable queries to find the referred object and text features will have fine-grained interaction with the visual features for reliable cross-modal reasoning.

## 3.2. Language as Queries

The key design comes from that we use a set of object queries conditioned on the language expression, termed *conditional queries*, as the Transformer decoder input. These queries are obligated to focus on the referred object only and produce the instance-aware dynamic kernels. The final segmentation masks are obtained by performing dynamic convolution between the dynamic kernels and their corresponding feature maps. Here, we adopt the Deformable-DETR [62] as our Transformer model due to its effectiveness and efficiency to capture the global pixel-level relations.

**Transformer Encoder.** First, a $1 \times 1$ convolution is applied on the multi-scale visual features $\mathcal{F}_v$ to reduce the channel dimension of all feature maps to $C = 256$. To enrich the information of visual features, we then incorporate projected visual features with the text feature $\mathcal{F}_e$ in a multiplication way and form the new multi-scale feature maps $\mathcal{F}_v' = \left\{ f_t' \right\}_{t=1}^T$. Afterwards, the fixed 2D positional encoding is added to feature maps of each frame and the summed features are fed into the Transformer encoder. To utilize the Transformer process the video frames independently, we flatten the spatial dimensions and move the temporal dimension to batch dimension for efficiency. Finally, the output of the Transformer encoder, *i.e.*, *encoded memory*, is then input to the decoder.

**Transformer Decoder.** We introduce $N$ object queries to represent the instances for each frame similar to [50], the difference lies in that the query weights are shared across video frames. This mechanism is more flexible to handle the length-variable videos and is more robust for the queries to track the same instances. Meanwhile, we repeat the sentence feature $f_e^s$ for $N$ times to fit the query number. Both



Figure 3. We visualize the predicted boxes from all the queries. It can be seen that the these boxes will locate near the referred object only even if there are other objects in the video.

the object queries and repeated sentence features are fed into the decoder as input. In this manner, all the queries will use the language expression as guidance and try to find the referred objects only. These *conditional queries* are duplicated to serve as the decoder input for all the frames and they are turned into instance embeddings by the decoder eventually, resulting in the set of $N_q = T \times N$ predictions. It should be noted the queries keep the same order across different frames and we refer to the queries in the same relative position (represented as the same shape in Figure 2) as *instance sequence* following [50]. Therefore, the temporal coherence of referred object could be achieved easily by linking the corresponding queries.

**Prediction Heads.** Three lightweight heads are built on top of the decoder to further transform the $N_q$ instance embeddings. The class head outputs the binary probability which indicates whether the instance is referred by the text sentence and this instance is visible in the current frame. It could also be modified to predict the referred object category by simply changing the output class number. The mask head is implemented by three consecutive linear layers. It produces the parameters of $N_q$ dynamic kernels $\Omega = \{\omega_i\}_{i=1}^{N_q}$, which is similar to the conditional convolutional filters in [45]. These parameters will be reshaped to form the three $1 \times 1$ convolution layers with the channel number as 8. The box head is a 3-layer feed forward network (FFN) with ReLU activation except for the last layer. It will predict the box location of the referred object and thus the position of dynamic kernels could be determined by the center of corresponding boxes.

**Dynamic Convolution.** Suppose now we have obtained the semantically-rich feature maps $\mathcal{F}_{\text{seg}} = \left\{ f_{seg}^t \right\}_{t=1}^T$ (will be discussed in Sec. 3.3) for each frame, the question is how we perform the instance sequence segmentation and obtain the masks of referred object from them. Since the dynamic kernels have captured the object-level information, we use them as convolution filters on the feature maps for instance decoding. Considering that the location prior of dynamic kernels $\Omega$ provides a strong and robust reference for the referred object, we concatenate the feature maps $\mathcal{F}_{seg}$ with relative coordinates for each dynamic kernel. Finally, the binary segmentation masks are generated by performing dynamic convolution between the conditional convolutional weights and their corresponding feature maps:
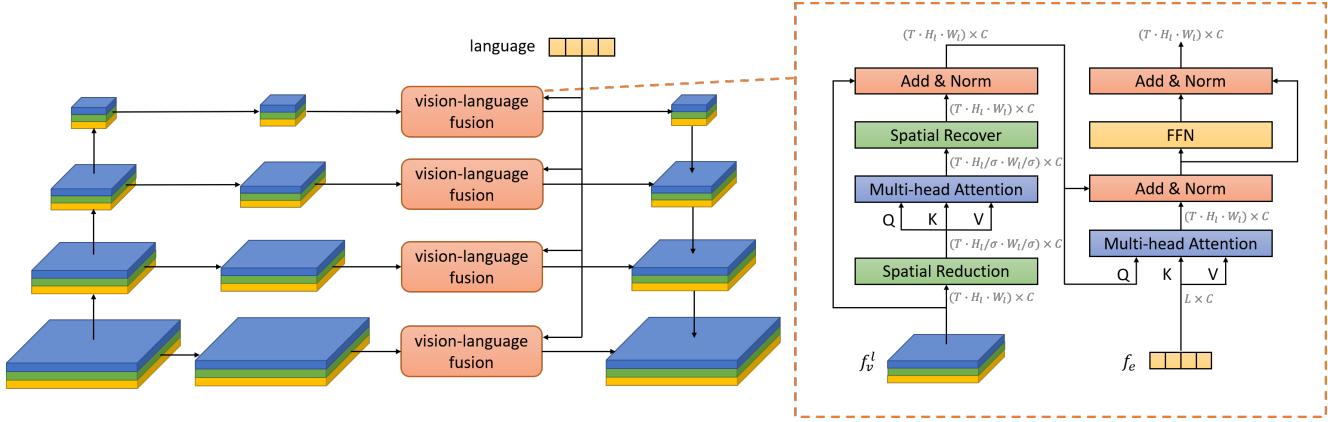
Figure 4. The architecture of cross-modal feature pyramid network (CM-FPN). Note that different colors in the feature maps represent different frames. The visual and textual features are interacted in all the levels of feature maps. The vision-language fusion process is illustrated in the dash box on the right.

$$\hat{s}_i = \left\{ \hat{f}_i \circledast \omega_i \right\}_{i=1}^{N_q} \tag{1}$$

where $\omega_i$ and $\hat{f}_i$ are the $i$-th dynamic kernel weights and its exclusive feature map, respectively. We reshape the output masks in frame-order sequence, resulting in a set as $\hat{\mathcal{S}} \in \mathbb{R}^{T \times N \times \frac{H}{4} \times \frac{W}{4}}$.

**Illustration of conditional queries.** It is well known that the decoder embedding and position embedding in Transformer decoder encode the content and spatial information respectively. In our framework, these two parts are fed with the text sentence feature and learanble queries parameters, so that all the queries are restricted by the language expression. As shown in figure 3, these queries will focus on the referred object only even if other objects exist in the video. And there will be one query with much higher score while the scores of other queries will be suppressed.

### 3.3. Cross-modal Feature Pyramid Network

Feature pyramid network (FPN) [25] is adopted to produce multi-scale feature maps for video frames. We construct a 4-level pyramid with the spatial stride from $4\times$ to $32\times$. Specifically, the first three stage features of Transformer encoded memory (with spatial strides $\{8, 16, 32\}$) and the $4\times$ feature from visual backbone are stacked to form the hierarchical features. Although the standard FPN can already provide a high-resolution feature map with rich visual semantics, such feature map lacks the linguistic information and would be sub-optimal for the cross-modal task. The previous work [42] only incorporates the language feature on the top level of FPN, which is a coarse fusion fashion. Here, we design a cross-modal feature pyramid network (CM-FPN) to perform multi-scale cross-modal fusion for finer interaction, as shown in Figure 4.

In each level, the interaction process is achieved by the vision-language fusion module. And we take the $l$-th level feature of FPN as an example to clarify the process. Here, we use $f_v^l \in \mathbb{R}^{T \times H_l \times W_l \times C}$ to represent the $l$-th level visual feature for simplicity. To model the spatio-temporal pixel-level relations of vision feature, we expect to feed it into a multi-head self-attention (MHSA) module. However, the computation of dense similarities make it intractable for the high-resolution feature maps. Inspired by [12, 49], we propose the *spatial reduction* and *spatial recover* operations to address the issue. Before the MHSA module, the spatial size of vision feature $F_v^l$ is downsampled by a factor of $\sigma$ while the temporal dimension is kept unchanged. Thus, the complexity of self-attention [46] operation would be greatly reduced, making the fusion module can be inserted into each level of FPN. Then, the spatial size of vision feature is recovered to $H_l \times W_l$ for maintaining fine-grained information. We set the downsample factors as $[8, 4, 2, 1]$ for the 4-level features maps. Next, $f_v^l$ interact with word-level feature $f_e$ in a *cross-attention* way, where the *query*, *key* are vision and language feature, respectively:

$$\text{Interact}(f_v^l, f_e) = \text{Softmax}\left( \frac{f_v^l W^Q \cdot (f_e W^K)^\mathrm{T}}{\sqrt{d_{\text{head}}}} \right) f_e W^V \tag{2}$$

where $W^Q, W^K, W^V \in \mathbb{R}^{C \times d_{\text{head}}}$ are learnable parameters. The visual feature plays the role of *query* in attention mechanism, and thus the pixels on the feature map that are strongly related to the language expression will be strengthened. We upsample and sum the cross-modal feature maps following the standard FPN top-down structure. Finally, we apply an additional $3 \times 3$ convolutional layer on the feature maps with spatial stride 4 to get the final feature maps $\mathcal{F}_{seg} = \left\{ f_{\text{seg}}^t \right\}_{t=1}^T$, where $f_{\text{seg}}^t \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C_d}$.

## 3.4. Instance Sequence Matching and Loss

Using $N$ *conditional queries*, we generate the set of $N_q = T \times N$ predictions, which can be regarded as the trajectories of $N$ instances on $T$ frames. As described previous, the predictions across frames maintain the same relative positions. Therefore, we can supervise the instance sequence as a whole using instance matching strategy [50]. Let us denote the prediction set as $\hat{y} = \{\hat{y}_i\}_{i=1}^N$, and the predictions for the $i$-th instance is represented by:

$$\hat{y}_i = \left\{\hat{p}_i^t, \hat{b}_i^t, \hat{s}_i^t\right\}_{t=1}^T \qquad (3)$$

For the $t$-th frame, $\hat{p}_i^t \in \mathbb{R}^1$ is a probability scalar indicating whether the instance corresponds to the referred object and this object is visible in the current frame. $\hat{b}_i^t \in \mathbb{R}^4$ is the normalized vector defines the center coordinates as well as the height and width of predicted box. $\hat{s}_i^t \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4}}$ is the predicted binary segmentation mask.

Since there is only one referred object in the video, the ground-truth instance sequence is represented as $y = \{c^t, b^t, s^t\}_{t=1}^T$. $c^t$ is an one-hot value and it equals 1 when the ground-truth instance is visible in the frame $I_t$ otherwise 0. To train the network, we first find the best prediction as the positive sample via minimizing the matching cost:

$$\hat{y}_{\text{pos}} = \arg\min_{\hat{y}_i \in \hat{y}} \mathcal{L}_{match}(y, \hat{y}_i) \qquad (4)$$

where

$$\mathcal{L}_{match}(y, \hat{y}_i) = \lambda_{cls}\mathcal{L}_{cls}(y, \hat{y}_i) + \lambda_{\text{box}}\mathcal{L}_{\text{box}}(y, \hat{y}_i) \\ + \lambda_{\text{mask}}\mathcal{L}_{\text{mask}}(y, \hat{y}_i) \qquad (5)$$

The matching cost is computed from each frame and normalized by the frame number. Here, $\mathcal{L}_{cls}(y, \hat{y}_i)$ is the focal loss [26] that supervises the predicted instance sequence reference results. The box-related loss sums up the L1 loss and GIoU loss [41]. And the mask-related loss is the combination of DICE loss [37] and binary mask focal loss. Both the two mask losses are spatio-temporally calculated over the entire video clip. The network is optimized by minimizing the total loss $\mathcal{L}_{match}$ for positive samples while letting the negative samples predict the $\varnothing$ class.

## 3.5. Inference

As mentioned previously, ReferFormer can handle the videos of arbitrary length in a single forward pass since all the frames share the same initial *conditional queries*. Given the video and language expression, ReferFormer will predict $N$ instance sequence. For each instance query, we average the predicted reference probabilities over all the frames and obtain the reference score set $\mathcal{P} = \{p_i\}_{i=1}^N$. We select the instance sequence with the highest average score and its index is denoted as $\sigma$:

$$\sigma = \arg\max_{i \in \{1,2,...,N\}} p_i \qquad (6)$$

The final segmentation masks for each frame $\mathcal{S} = \{s_t\}_{t=1}^T$ is obtained from the mask candidates set $\hat{\mathcal{S}}$ by selecting the corresponding queries indexed with $\sigma$. No post-process is needed for associating objects since the linked queries naturally track the same instance.

# 4. Experiments

## 4.1. Datasets and Metrics

**Datasets.** The experiments are conducted on the four popular R-VOS benchmarks: Ref-Youtube-VOS [42], Ref-DAVIS17 [20], A2D-Sentences and JHMDB-Sentences [13]. Ref-Youtube-VOS [42] is a large-scale benchmark which covers 3,978 videos with ~15K language descriptions. Ref-DAVIS17 [20] is built upon DAVIS17 [40] by providing the language description for a specific object in each video and contains 90 videos. A2D-Sentences and JHMDB-Sentences are created by providing the additional textual annotations on the original A2D [52] and JHMDB [18] datasets. A2D-Sentences contains 3,782 videos and each video has 3-5 frames annotated with the pixel-level segmentation masks. JHMDB-Sentences has 928 videos with the 928 corresponding sentences in total.

**Evaluation Metrics.** We use the standard evaluation metrics for Ref-Youtube-VOS and Ref-DAVIS17: region similarity ($\mathcal{J}$), contour accuracy ($\mathcal{F}$) and their average value ($\mathcal{J}\&\mathcal{F}$). For Ref-Youtube-VOS, as the annotations of validation set are not released publicly, we evaluate our method on the official challenge server [1]. Ref-DAVIS17 is evaluated by the official evaluation code [2].

On A2D-Sentences and JHMDB-Sentences, the model is evaluated with the criteria of Precision@K, Ovrall IoU, Mean IoU and mAP over 0.50:0.05:0.95. The Precision@K measures the percentage of test samples whole IoU scores are higher than the threshold K. Following standard protocol, the thresholds are set as 0.5:0.1:0.9.

## 4.2. Implementation Details.

**Model Settings.** We test our models under different visual backbones including: ResNet [14], Swin Transformer [30] and Video Swin Transformer [31]. The text encoder is selected as RoBERTa [29] and its parameters are frozen during the entire training stage. Following [62], we use the last stage features from the visual backbone as the input to Transformer, their corresponding spatial strides are $\{8, 16, 32\}$. In the Transformer model, we adopt 4 encoder

| Method | Backbone | Ref-Youtube-VOS | | | Ref-DAVIS17 | | |
|---|---|---|---|---|---|---|---|
| | | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
| **Spatial Visual Backbones** | | | | | | | |
| CMSA [57] | ResNet-50 | 34.9 | 33.3 | 36.5 | 34.7 | 32.2 | 37.2 |
| CMSA + RNN [57] | ResNet-50 | 36.4 | 34.8 | 38.1 | 40.2 | 36.9 | 43.5 |
| URVOS [42] | ResNet-50 | 47.2 | 45.3 | 49.2 | 51.5 | 47.3 | 56.0 |
| ReferFormer | ResNet-50 | 55.6 | 54.8 | 56.5 | **58.5** | **55.8** | **61.3** |
| ReferFormer* | ResNet-50 | **58.7** | **57.4** | **60.1** | - | - | - |
| PMINet [10] | ResNeSt-101 | 48.2 | 46.7 | 49.6 | - | - | - |
| PMINet + CFBI [10] | ResNeSt-101 | 53.0 | 51.5 | 54.5 | - | - | - |
| CITD* [23] | ResNet-101 | 56.4 | 54.8 | 58.1 | - | - | - |
| ReferFormer | ResNet-101 | 57.3 | 56.1 | 58.4 | - | - | - |
| ReferFormer* | ResNet-101 | **59.3** | **58.1** | **60.4** | - | - | - |
| PMINet + CFBI [10] | Ensemble | 54.2 | 53.0 | 55.5 | - | - | - |
| CITD [23] | Ensemble | 61.4 | 60.0 | 62.7 | - | - | - |
| ReferFormer | Swin-L | 62.4 | 60.8 | 64.0 | 60.5 | 57.6 | 63.4 |
| ReferFormer* | Swin-L | **64.2** | **62.3** | **66.2** | - | - | - |
| **Spatio-temporal Visual Backbones** | | | | | | | |
| MTTR$^{\dagger}$ ($\omega = 12$) [3] | Video-Swin-T | 55.3 | 54.0 | 56.6 | - | - | - |
| ReferFormer $^{\dagger}$ ($\omega = 5$) | Video-Swin-T | **56.0** | **54.8** | **57.3** | - | - | - |
| ReferFormer | Video-Swin-T | 59.4 | 58.0 | 60.9 | - | - | - |
| ReferFormer* | Video-Swin-T | **62.6** | **59.9** | **63.3** | - | - | - |
| ReferFormer | Video-Swin-S | 60.1 | 58.6 | 61.6 | - | - | - |
| ReferFormer* | Video-Swin-S | **63.3** | **61.4** | **65.2** | - | - | - |
| ReferFormer | Video-Swin-B | 62.9 | 61.3 | 64.6 | 61.1 | 58.1 | 64.1 |
| ReferFormer* | Video-Swin-B | **64.9** | **62.8** | **67.0** | - | - | - |

Table 1. Comparison with the state-of-the-art methods on Ref-Youtube-VOS and Ref-DAVIS17. * means joint trainig with Ref-COCO dataset. † indicates the spatio-temporal visual backbone is trained from scratch.

layers and 4 decoder layers and the hidden dimension is $C = 256$. The number of *conditional query* is set as 5 otherwise specified.

**Training Details.** During training, we use sliding-windows to obtain the clips from a video and each clip consist of 5 randomly sampled frames. Following [50], the data augmentation includes random horizontal flip, random resize, random crop and photometric distortion. All frames are downsampled so that the short side has the size of 360 and the maximum size for the long side is 640 to fit GPU memory. The coefficients for losses are set as $\lambda_{cls} = 2$, $\lambda_{L1} = 5$, $\lambda_{giou} = 2$, $\lambda_{dice} = 5$, $\lambda_{focal} = 2$.

Most of our experiments follow the pretrain-then-finetune process. And some models are trained from scratch for fair comparison. Additionally, on Ref-Youtube-VOS, we also reports the results by training the mixed data from Ref-Youtube-VOS and Ref-COCO [59]. The joint training technique has proven the effectiveness in many VIS tasks [1, 24, 51]. Please see more in the supplementary materials.

**Inference Details.** During inference, the video frames are downscaled to 360p. We directly output the predicted seg-mentation masks without post-process. On Ref-Youtube-VOS, we further use a simple post-process technique to refine the object masks. Concretely, we first select a frame with the highest prediction score as the reference frame. Then, we apply the off-the-shelf mask propagation method CFBI [56] to propagate the predicted mask of this frame forward and backward to the entire video. The results with post-process are shown in Table 5.

### 4.3. Main Results

**Ref-Youtube-VOS & Ref-DAVIS17** We compare our method with other state-of-the-art methods in Table 1. CITD [23] and PMINet [10] are the top-2 solutions in 2021 Ref-Youtube-VOS Challenge. Their ensemble results are based on building 5 and 4 models, respectively. It can be observed that ReferFormer outperforms previous methods on the two datasets under all metrics and with a large marge. On Ref-Youtube-VOS, ReferFormer with a ResNet-50 backbone achieves the overall $\mathcal{J}\&\mathcal{F}$ of 55.6, which is 8.4 points higher than the previous state-of-the-art work URVOS [42], and even beats PMINet [10] using the ensemble models and adopting post-process (55.6 vs

| Method | Backbone | Precision | | | | | IoU | | mAP |
|---|---|---|---|---|---|---|---|---|---|
| | | P@0.5 | P@0.6 | P@0.7 | P@0.8 | P@0.9 | Overall | Mean | |
| Hu *et al.* [15] | VGG-16 | 34.8 | 23.6 | 13.3 | 3.3 | 0.1 | 47.4 | 35.0 | 13.2 |
| Gavrilyuk *et al.* [13] | I3D | 47.5 | 34.7 | 21.1 | 8.0 | 0.2 | 53.6 | 42.1 | 19.8 |
| CMSA + CFSA [58] | ResNet-101 | 48.7 | 43.1 | 35.8 | 23.1 | 5.2 | 61.8 | 43.2 | - |
| ACAN [48] | I3D | 55.7 | 45.9 | 31.9 | 16.0 | 2.0 | 60.1 | 49.0 | 27.4 |
| CMPC-V [28] | I3D | 65.5 | 59.2 | 50.6 | 34.2 | 9.8 | 65.3 | 57.3 | 40.4 |
| ClawCraneNet [22] | ResNet-50/101 | 70.4 | 67.7 | 61.7 | 48.9 | 17.1 | 63.1 | 59.9 | - |
| MTTR ($\omega = 8$) [3] | Video-Swin-T | 72.1 | 68.4 | 60.7 | 45.6 | 16.4 | 70.2 | 61.8 | 44.7 |
| MTTR ($\omega = 10$) [3] | Video-Swin-T | 75.4 | 71.2 | 63.8 | 48.5 | 16.9 | 72.0 | 64.0 | 46.1 |
| ReferFormer$^\dagger$ ($\omega = 6$) | Video-Swin-T | 76.0 | 72.2 | 65.4 | 49.8 | 17.9 | 72.3 | 64.1 | 48.6 |
| ReferFormer ($\omega = 5$) | Video-Swin-T | 82.8 | 79.2 | 72.3 | 55.3 | 19.3 | 77.6 | 69.6 | 52.8 |
| ReferFormer ($\omega = 5$) | Video-Swin-S | 82.6 | 79.4 | 73.1 | 57.4 | 21.1 | 77.7 | 69.8 | 53.9 |
| **ReferFormer** ($\omega = 5$) | Video-Swin-B | **83.1** | **80.4** | **74.1** | **57.9** | **21.2** | **78.6** | **70.3** | **55.0** |

Table 2. Comparison with the state-of-the-art methods on A2D-Sentences. $^\dagger$ means our model is trained from scratch.

| Method | Backbone | Precision | | | | | IoU | | mAP |
|---|---|---|---|---|---|---|---|---|---|
| | | P@0.5 | P@0.6 | P@0.7 | P@0.8 | P@0.9 | Overall | Mean | |
| Hu *et al.* [15] | VGG-16 | 63.3 | 35.0 | 8.5 | 0.2 | 0.0 | 54.6 | 52.8 | 17.8 |
| Gavrilyuk *et al.* [13] | I3D | 69.9 | 46.0 | 17.3 | 1.4 | 0.0 | 54.1 | 54.2 | 23.3 |
| CMSA + CFSA [58] | ResNet-101 | 76.4 | 62.5 | 38.9 | 9.0 | 0.1 | 62.8 | 58.1 | - |
| ACAN [48] | I3D | 75.6 | 56.4 | 28.7 | 3.4 | 0.0 | 57.6 | 58.4 | 28.9 |
| CMPC-V [28] | I3D | 81.3 | 65.7 | 37.1 | 7.0 | 0.0 | 61.6 | 61.7 | 34.2 |
| ClawCraneNet [22] | ResNet-50/101 | 88.0 | 79.6 | 56.6 | 14.7 | 0.2 | 64.4 | 65.6 | - |
| MTTR ($\omega = 8$) [3] | Video-Swin-T | 91.0 | 81.5 | 57.0 | 14.4 | 0.1 | 67.4 | 67.9 | 36.6 |
| MTTR ($\omega = 10$) [3] | Video-Swin-T | 93.9 | 85.2 | 61.6 | 16.6 | 0.1 | 70.1 | 69.8 | 39.2 |
| ReferFormer$^\dagger$ ($\omega = 6$) | Video-Swin-T | 93.3 | 84.2 | 61.4 | 16.4 | 0.3 | 70.0 | 69.3 | 39.1 |
| ReferFormer ($\omega = 5$) | Video-Swin-T | 95.8 | 89.3 | 66.8 | 18.9 | 0.2 | 71.9 | 71.0 | 42.2 |
| ReferFormer ($\omega = 5$) | Video-Swin-S | 95.8 | 90.1 | 68.7 | 20.3 | 0.2 | 72.8 | 71.5 | 42.4 |
| **ReferFormer** ($\omega = 5$) | Video-Swin-B | **96.2** | **90.2** | **70.2** | **21.0** | **0.3** | **73.0** | **71.8** | **43.7** |

Table 3. Comparison with the state-of-the-art methods on JHMDB-Sentences. $^\dagger$ means our model is trained from scratch.

54.2). Using the strong Swin-Large [30] backbone, Refer-Former reaches the surprising 62.4 $\mathcal{J}\&\mathcal{F}$ without bells and whistles, which obviously exceeds the ensemble results of the complicated, multi-stage method CITD [23]. By using the joint training process, the performance of our model can be further boosted to 64.2 $\mathcal{J}\&\mathcal{F}$, creating a fairly high new record. Additionally, we also test the Video Swin Transformer [31] as the backbones. It is well known that the spatio-temporal visual encoder has strong ability to capture both the spatial characteristics and the temporal cues. For a fair comparison with MTTR [3], we train our model with the Video-Swin-Tiny backbone from scratch. It can be seen that our method outperforms MTTR under all the metrics with the smaller window size (5 vs 12). Comparing the results of ReferFormer under Video-Swin-Tiny backbone, it proves that the model benefits from the pretraining stage and joint training process to address the overfitting issue.

On Ref-DAVIS17, our method also achieves the best results under the same ResNet-50 setting (58.5 $\mathcal{J}\&\mathcal{F}$). And the performance consistently improves by using stronger backbones, which proves the generality of our method.

**A2D-Sentences & JHMDB-Sentences** We further evaluate our method on the A2D-Sentences dataset and compare the performance with other state-of-the-art methods in Table 2. ClawCraneNet [22] is a mutli-stage method which use the off-the-shelf instance segmentation model (with ResNet-101 backbone) to provide the mask candidates. From Table 2, it is obvious that our method achieves the impressive improvement over the previous methods. Compared with the recent MTTR [3], our method exhibits the clear performance advantage (+2.5 mAP) with smaller window size (6 vs. 10). Incorporating the pretraining stage, ReferFormer with Video-Swin-Base visual backbone achieves 55.0 mAP which shows a significant gain of 8.9 mAP over previous best result. And ReferFormer also demonstrates its strong ability to produce high-quality masks via the stringent metrics (*e.g.*, 57.9 for P@0.8 and 21.2 for P@0.9).

| Components | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|
| Baseline | 47.2 (−7.6) | 50.1 (−7.2) |
| w/o Visual-language Fusion | 53.0 (−1.8) | 56.2 (−1.1) |
| w/o Relative Coordinates | 53.7 (−1.1) | 55.9 (−1.4) |
| Full Model | 54.8 | 57.3 |

Table 4. Ablation study on the components of ReferFormer. The visual backbone is Video-Swin-Tiny.

| Backbone | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|
| ResNet-50 | 55.6 | 54.8 | 56.5 |
| ResNet-50* | 59.4 (+3.8) | 58.1 (+3.3) | 60.8 (+4.3) |
| ResNet-101 | 57.3 | 56.1 | 58.4 |
| ResNet-101* | 60.3 (+3.0) | 58.8 (+2.7) | 61.8 (+3.4) |
| Swin-T | 58.7 | 57.6 | 59.9 |
| Swin-T* | 61.2 (+2.5) | 59.7 (+2.1) | 62.6 (+2.7) |
| Swin-S | 59.6 | 58.1 | 61.1 |
| Swin-S* | 61.3 (+1.7) | 59.7 (+1.6) | 63.0 (+1.9) |
| Swin-B | 61.8 | 60.1 | 63.4 |
| Swin-B* | 63.1 (+1.3) | 61.4 (+1.3) | 64.8 (+1.4) |
| Swin-L | 62.4 | 60.8 | 64.0 |
| Swin-L* | 63.3 (+0.9) | 61.6 (+0.8) | 65.1 (+1.1) |

Table 5. Ablation study on the visual backbones. * indicates using CFBI [56] as post-process.

We also evaluate the models on JHMDB-Sentences without finetuning to further prove the generality of our method. As shown in Table 3, ReferFormer significantly outperforms all the existing methods. It is noticeable that all the methods produce low scores on P@0.9. A possible reason is that the ground-truth masks are generated from human puppets, leading to the inaccurate mask annotations.

## 4.4. Ablation Study

In this section, we perform extensive ablation studies on Ref-Youtube-VOS to study the effect of core components in our model. All models are based on Video-Swin-Tiny visual backbone and we train the models from scratch otherwise specified. The detailed analysis is as follows.

**Component Analysis.** We build a simple Transformer bottom-up baseline. Specifically, considering a video clip of $T$ frames, we flatten the temporal and spatial dimension into one dimension and then concatenate the visual features with the textual features along length dimension to form the multi-modal feature map $f_m \in \mathbb{R}^{(T \times H \times W + L) \times C}$. The vanilla Transformer encoder builds the global dependencies between the visual and textual features. Afterwards, we extract the visual features from the encoded memory and construct a standard FPN-like architecture upon them for generating the segmentation masks. The baseline method operates the fixed-length video of 5 frames during the training and inference phases. We report the performance of the baseline method and also study the effect of core components in Table 4.

First, from the first row of Table 4, the baseline method only achieves 47.2 $\mathcal{J}$ and 50.1 $\mathcal{F}$. This inferior behavior attributes to two reasons: (1) The baseline method can not distinguish the similar objects that are close together and tends to segment the most salient region. In contrast, our method performs well with only 1 conditional query (see Table 6(a)), proving that dynamic convolution is essential for segmenting the referred object. (2) Our method uses a set of shared queries to track instances in all frames, and the best query is determined by the *voting* scores of each frame. In this sense, our model can produce a reliable reasoning result and keep the temporal consistency in the entire video. On the contrary, the baseline method could be regarded as a image-level method that independently predicts the results of each frame even though the model is able to aggregate the information from other frames.

Second, comparing the second and last row of Table 4, we can see that the standard FPN has already achieved strong performance and the vision-language fusion process further helps to provide more accurate segmentation. This is because the object mask would be inaccurate due to light variation, whereas the cross-modal fusion uses the text as a complementary to strengthen the object pixel features and thus facilitates the segmentation prediction. Another technique is concatenating the relative coordinates of dynamic kernels with the mask features, this would help the model better determine the location of referred object and lead to performance improvement, as shown in the third row in Table 4.

**Visual Backbone.** We implement different visual backbones and report the results in Table 5. As expected, the performance of our model consistently increases by using stronger backbones. And the CFBI [56] post-process can help to further boost the performance under all backbone settings. Interestingly, we observe that the performance improvement by post-process tends to narrow when the backbone gets stronger, *e.g.*, +3.8 for ResNet-50 and +0.9 for Swin-Large when considering the $\mathcal{J}\&\mathcal{F}$ metric. This phenomenon shows that the visual encoder is essential for providing reliable reasoning on which object is described and generating the precise masks.

**Number of Conditional Queries.** Benefit from the design of *conditional queries*, all the initial object queries tend to find the referred object only. In this situation, we can only use a relatively small number of queries. In Table 6(a), we study the effect of query number for each frame. It can be seen that the model achieves considerable results under all these settings, even with $N = 1$. Certainly, more queries enable the model make judgement from a wide range of instance candidates, which could better handle the com-

(a)

a man riding a motorbike    a green motorbike

(b)

a person is showing his skate board skills on the road    a skate board carrying a person skating on the road

Figure 5. Visualization results on (a) Ref-DAVIS17 and (b) Ref-Youtube-VOS. Our unified framework is able to detect, segment and track the referred object simultaneously.

| Queries | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|
| 1 | 53.6 | 52.7 | 54.5 |
| 3 | 54.2 | 53.2 | 55.2 |
| 5 | 56.0 | 54.8 | 57.3 |
| 8 | 55.3 | 54.1 | 56.6 |

(a) The effect of query number.

| Frames | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|
| 1 | 50.0 | 48.4 | 51.6 |
| 3 | 54.8 | 53.6 | 56.0 |
| 5 | 56.0 | 54.8 | 57.3 |

(b) The effect of frame number.

| Class | Box | Mask | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|---|---|
| ✓ | ✓ | | 55.2 | 54.0 | 56.4 |
| ✓ | | ✓ | 54.5 | 53.5 | 55.5 |
| ✓ | ✓ | ✓ | 56.0 | 54.8 | 57.3 |

(c) The effect of label assignment method.

Table 6. Ablation study on different settings of ReferFormer. All the models are using Video-Swin-Tiny as visual backbone.

plicated scenes where the similar objects are clustered together. The performance saturates at $N = 5$ and begins to slightly decrease when the query number gets larger. We conjecture that it is caused by the imbalance of label assignment as there is only one positive sample in each frame.

**Number of Training Clip Frames.** We study the effect of training clip frame number in Table 6(b). Note that under $T = 1$, the model can be viewed as an image-level method and the performance of metric $\mathcal{J}\&\mathcal{F}$ is only 50.0. When the frame number increases to 3, the model enjoys an significant $\mathcal{J}\&\mathcal{F}$ gain of 4.8. This is because using more frames to form a clip helps the model better aggregate the temporal action-related information. We choose $T = 5$ by default.

**Label Assignment Method.** Our framework is able to predict the reference probability, box location and segmentation mask of the referred object. We find the optimal positive sample by minimizing the overall matching cost in Eq.4. There are some variants in the label assignment

method and we carry out the comparison experiments in Table 6(c). From the first two rows in Table 6(c) we show that the lack of box or mask cost would both lead to the performance drop. With the segmentation-centric design, the mask cost is the most direct guidance for optimization, and the box provides the location prior for dynamic kernel. Thus, the combination of classification, box and mask cost shows more robustness.

### 4.5. Visualization Results

We show the visualization results of our model in Figure 5. It can be seen that ReferFormer is able to segment and track the referred object in challenging cases, *e.g.*, person pose variations, instances occlusion and instances that are partially displayed or completely disappeared in the camera.

## 5. Conclusion

In this work, we propose ReferFormer, an extremely simple and unified framework for referring video object

segmentation. This framework provides a new perspective for the R-VOS task which views the language as queries. These queries are restricted to attend to the referred object only, and the object tracking is easily achieved by linking the corresponding queries. Given the video clip and an expression, our framework directly produces the segmentation masks as well as the detected boxes of the referred object in all frames without post-process. We validate our model on Ref-Youtube-VOS, Ref-DAVIS17, A2D-Sentences and JHMDB-Sentences and it shows the state-of-the-art performance on the four benchmarks.

# References

[1] Ali Athar, Sabarinath Mahadevan, Aljosa Osep, Laura Leal-Taixé, and Bastian Leibe. Stem-seg: Spatio-temporal embeddings for instance segmentation in videos. In *European Conference on Computer Vision*, pages 158–177. Springer, 2020. 7

[2] Miriam Bellver, Carles Ventura, Carina Silberer, Ioannis Kazakos, Jordi Torres, and Xavier Giro-i Nieto. Refvos: A closer look at referring expressions for video object segmentation. *arXiv preprint arXiv:2010.00263*, 2020. 2

[3] Adam Botach, Evgenii Zheltonozhskii, and Chaim Baskin. End-to-end referring video object segmentation with multimodal transformers. *arXiv preprint arXiv:2111.14821*, 2021. 3, 7, 8

[4] 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. 2, 3

[5] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4974–4983, 2019. 2

[6] Xin Chen, Bin Yan, Jiawen Zhu, Dong Wang, Xiaoyun Yang, and Huchuan Lu. Transformer tracking. In *CVPR*, 2021. 3

[7] 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

[8] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. *arXiv preprint arXiv:2106.05210*, 2021. 2

[9] Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. Vision-language transformer and query generation for referring segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16321–16330, 2021. 2

[10] Zihan Ding, Tianrui Hui, Shaofei Huang, Si Liu, Xuan Luo, Junshi Huang, and Xiaoming Wei. Progressive multimodal interaction network for referring video object segmentation. 2021. 7

[11] 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. *arXiv preprint arXiv:2010.11929*, 2020. 3

[12] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. *arXiv preprint arXiv:2104.11227*, 2021. 5

[13] Kirill Gavrilyuk, 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. 6, 8

[14] 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. 4, 6

[15] Ronghang Hu, Marcus Rohrbach, and Trevor Darrell. Segmentation from natural language expressions. In *European Conference on Computer Vision*, pages 108–124. Springer, 2016. 8

[16] Shaofei Huang, Tianrui Hui, Si Liu, Guanbin Li, Yunchao Wei, Jizhong Han, Luoqi Liu, and Bo Li. Referring image segmentation via cross-modal progressive comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10488–10497, 2020. 2

[17] Sukjun Hwang, Miran Heo, Seoung Wug Oh, and Seon Joo Kim. Video instance segmentation using inter-frame communication transformers. *arXiv preprint arXiv:2106.03299*, 2021. 2, 3

[18] Hueihan Jhuang, Juergen Gall, Silvia Zuffi, Cordelia Schmid, and Michael J Black. Towards understanding action recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 3192–3199, 2013. 6

[19] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790, 2021. 3

[20] 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. 6, 14

[21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25:1097–1105, 2012. 2

[22] Chen Liang, Yu Wu, Yawei Luo, and Yi Yang. Clawcranenet: Leveraging object-level relation for text-based video segmentation. *arXiv preprint arXiv:2103.10702*, 2021. 8

[23] Chen Liang, Yu Wu, Tianfei Zhou, Wenguan Wang, Zongxin Yang, Yunchao Wei, and Yi Yang. Rethinking cross-modal interaction from a top-down perspective for referring video object segmentation. *arXiv preprint arXiv:2106.01061*, 2021. 2, 3, 7, 8

[24] Huaijia Lin, Ruizheng Wu, Shu Liu, Jiangbo Lu, and Ji-aya Jia. Video instance segmentation with a propose-reduce paradigm. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1739–1748, 2021. 7

[25] 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

[26] 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. 6

[27] 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 *European conference on computer vision*, pages 740–755. Springer, 2014. 2

[28] Si Liu, Tianrui Hui, Shaofei Huang, Yunchao Wei, Bo Li, and Guanbin Li. Cross-modal progressive comprehension for referring segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 8

[29] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettle-moyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 4, 6

[30] 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. *arXiv preprint arXiv:2103.14030*, 2021. 3, 6, 8

[31] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. *arXiv preprint arXiv:2106.13230*, 2021. 4, 6, 8

[32] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 1

[33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 14

[34] Gen Luo, Yiyi Zhou, Xiaoshuai Sun, Liujuan Cao, Chenglin Wu, Cheng Deng, and Rongrong Ji. Multi-task collaborative network for joint referring expression comprehension and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10034–10043, 2020. 2

[35] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016. 14

[36] Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixe, and Christoph Feichtenhofer. TrackFormer: Multi-object tracking with transformers. *arXiv preprint arXiv:2101.02702*, 2021. 3

[37] 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. 6

[38] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9226–9235, 2019. 2

[39] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732, 2016. 1

[40] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017. 6, 14

[41] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 658–666, 2019. 6

[42] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object segmentation network with a large-scale benchmark. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, pages 208–223. Springer, 2020. 3, 5, 6, 7, 14

[43] Peize Sun, Yi Jiang, Rufeng Zhang, Enze Xie, Jinkun Cao, Xinting Hu, Tao Kong, Zehuan Yuan, Changhu Wang, and Ping Luo. Transtrack: Multiple-object tracking with transformer. *arXiv preprint arXiv:2012.15460*, 2020. 3

[44] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chen-feng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14454–14463, 2021. 3

[45] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 282–298. Springer, 2020. 4

[46] 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. 2, 3, 5

[47] Paul Voigtlaender, Yuning Chai, Florian Schroff, Hartwig Adam, Bastian Leibe, and Liang-Chieh Chen. Feelvos: Fast end-to-end embedding learning for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9481–9490, 2019. 2

[48] Hao Wang, Cheng Deng, Junchi Yan, and Dacheng Tao. Asymmetric cross-guided attention network for actor and action video segmentation from natural language query. In

*Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3939–3948, 2019. 8

[49] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. *arXiv preprint arXiv:2102.12122*, 2021. 5

[50] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8741–8750, 2021. 2, 3, 4, 6, 7

[51] Junfeng Wu, Yi Jiang, Wenqing Zhang, Xiang Bai, and Song Bai. Seqformer: a frustratingly simple model for video instance segmentation. *arXiv preprint arXiv:2112.08275*, 2021. 2, 3, 7

[52] 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. 6

[53] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018. 1

[54] Bin Yan, Houwen Peng, Jianlong Fu, Dong Wang, and Huchuan Lu. Learning spatio-temporal transformer for visual tracking. In *ICCV*, 2021. 3

[55] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5188–5197, 2019. 3

[56] Zongxin Yang, Yunchao Wei, and Yi Yang. Collaborative video object segmentation by foreground-background integration. In *European Conference on Computer Vision*, pages 332–348. Springer, 2020. 2, 7, 9

[57] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. Cross-modal self-attention network for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10502–10511, 2019. 7

[58] Linwei Ye, Mrigank Rochan, Zhi Liu, Xiaoqin Zhang, and Yang Wang. Referring segmentation in images and videos with cross-modal self-attention network. *arXiv preprint arXiv:2102.04762*, 2021. 8

[59] 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. 2, 7, 14

[60] Yuting Zhang, Luyao Yuan, Yijie Guo, Zhiyuan He, I-An Huang, and Honglak Lee. Discriminative bimodal networks for visual localization and detection with natural language queries. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 557–566, 2017. 2

[61] 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

[62] 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. 2, 3, 4, 6

## A. Additional Dataset Details

**Ref-Youtube-VOS** [42] is a large-scale benchmark which covers 3,978 videos with ∼15K language descriptions. There are 3,471 videos with 12,913 expressions in training set and 507 videos with 2,096 expressions in validation set. According to the R-VOS competition, videos in the validation set are further split into 202 and 305 videos for the competition validation and test purpose. Since the test server is currently inaccessible, the results are reported by submitting our predictions to the validation server[3].

**Ref-DAVIS17** [20] is built upon DAVIS17 [40] by providing the language description for a specific object in each video. It contains 90 videos with 1,544 expression sentences describing 205 objects in total. The dataset is split into 60 videos and 30 videos for training and validation, respectively. Since there are two annotators and each of them gives the *first-frame* and *full-video* textual description for one referred object, we report the results by averaging the scores using the official evaluation code [4].

## B. Additional Implementation Details

Our model is optimized using AdamW [33] optimizer with the weight decay of $5 \times 10^{-4}$, initial learning rate of $5 \times 10^{-5}$ for visual backbone and $10^{-4}$ for the rest. We first pretrain our model on the image referring segmentation datasets Ref-COCO [59], Ref-COCOg [59] and Ref-COCO+ [35] by setting $T = 1$ with the batch size of 2 on each GPU. The pretrain procedure runs for 12 epochs with the learning rate decays divided by 10 at epoch 8 and 10. Then, on Ref-Youtube-VOS, we finetune the model for 6 epochs with 1 video clip per GPU. The learning rate decays by 10 at the 3-th and 5-th epoch. On Ref-DAVIS17, we directly report the results using the model trained on Ref-Youtube-VOS without finetune.

For A2D-Sentences, we feed the model with the window size of 5. The model is finetuned for 6 epochs with the learning rate decays at the 3-th and 5-th epoch by a factor of 0.1. On JHMDB-Sentences, following the previous works, we evaluate the generality of our method using the model trained on A2D-Sentences without finetune.

Additionally, on the Ref-Youtube-VOS, we also adopt the joint training technique by mixing the dataset with Ref-COCO/+/g. Specifically, for each image in the Ref-COCO dataset, we augment it with $\pm 20°$ to form a 5-frame pseudo video clip. The joint training takes 12 epochs with the learning rate decays at the 8-th and 10-th epoch by a factor of 0.1. We use 32 V100 GPUS for the joint training and each GPU is fed with 2 video clips. It should be noted that the text encoder is froze all the time.

---

## C. Additional Details of Dynamic Convolution

We give the pseudo-code of dynamic convolution in Figure C1, where we take one dynamic kernel for clarification. Specifically, a linear projection is applied to transform the instance embedding into dynamic convolutional weights. Then, the mask features pass through consecutive dynamic convolutional layers with the ReLU activation function. There is no normalization or activation after the last dynamic convolutional layer, and the output channel number of last layer is 1.

```python
def dynamic_convolution(mask_feats, dynamic_feats):
    # mask_feats: (B, C, H/4, W/4)
    # dynamic_feats: (B, C)

    # parameters of dynamic convs: (B, num_params)
    dynamic_params = linear(dynamic_features)
    # parse conv parameters
    # weights[l]: (out_channels, in_channels, 1, 1)
    # bias[l]: (out_channels)
    weights, bias = parse_dynamic_params(dynaic_params)

    # dynamic convolution
    n_layer = len(weights)
    x = mask_feats
    for i, (x, b) in enumerate(zip(weights, bias)):
        x = conv2d(x, w, bias=b, stride=1, padding=0)
        if i < n_layer - 1:
            x = relu(x)
    # x: (B, H/4, W/4)
    return x
```

Figure C1. Pseudo-code of dynamic convolution, we take one dynamic kernel for clarification. For multiple dynamic kernels, we use group convolution in `conv2d` for efficient implementation. `linear`: linear projection.

## D. Additional Experiment Results

By default, our models are trained in the class-agnostic way, *i.e.*, decide whether the object is referred or not. As described in Sec 3.2, the class head can be easily modified to predict the referred object category by simply change the class number. In this way, we train our model in a class-discriminative way and show the results in Table D1. We could observe the class-agnostic training method has clear performance gain ($+2.1 \mathcal{J}\&\mathcal{F}$) over the strong class-discriminative training results, since the binary classification is easier to optimize. The selection of training method can flexibly depend on the usage in real applications.

| Class Agnostic | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|
|  | 53.9 | 52.8 | 55.0 |
| ✓ | 56.0 | 54.8 | 57.3 |

Table D1. Ablation study on the class-agnostic training.