
R²-VOS: Robust Referring Video Object Segmentation via Relational Cycle Consistency

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Referring video object segmentation (R-VOS) aims to segment the object masks in a
2 video given a referring linguistic expression to the object. It is a recently introduced
3 task attracting growing research attention. However, all existing works make a
4 strong assumption: The object depicted by the expression must exist in the video,
5 namely, the expression and video must have an object-level semantic consensus.
6 This is often violated in real-world applications where an expression can be queried
7 to false videos, and existing methods always fail in such false queries due to abusing
8 the assumption. In this work, we emphasize that studying semantic consensus is
9 necessary to improve the robustness of R-VOS. Accordingly, we pose an extended
10 task from R-VOS without the semantic consensus assumption, named Robust R-
11 VOS (R²-VOS). The R²-VOS task is essentially related to the joint modeling of the
12 primary R-VOS task and its dual problem (text reconstruction). We embrace the
13 observation that the embedding spaces have relational consistency through the cycle
14 of text-video-text transformation, which connects the primary and dual problems.
15 We leverage the cycle consistency to discriminate the semantic consensus, thus
16 advancing the primary task. Parallel optimization of the primary and dual problems
17 are enabled by introducing an early grounding medium. A new evaluation dataset,
18 R²-Youtube-VOS, is collected to measure the robustness of R-VOS models against
19 unpaired videos and expressions. Extensive experiments demonstrate that our
20 method not only identifies negative pairs of unrelated expressions and videos,
21 but also improves the segmentation accuracy for positive pairs with a superior
22 disambiguating ability. Our model achieves the state-of-the-art performance on
23 Ref-DAVIS17, Ref-Youtube-VOS, and the novel R²-Youtube-VOS dataset.

24 1 Introduction

25 Referring video object segmentation (R-VOS) aims to segment a referred object in a video sequence
26 given a linguistic expression. R-VOS has witnessed growing interest thanks to its promising potential
27 in human-computer interaction applications such as video editing and augmented reality. Unlike
28 other video segmentation tasks [45, 36, 35, 46] that only rely on visual cues, R-VOS [13] pairs a
29 target video with a linguistic expression referring to an object.

30 Previous works [1, 44] tackle the R-VOS problem with a strong assumption that the referred object
31 exists in the video, i.e., there is an object-level semantic consensus between the expression and the
32 video. However, this assumption does not always hold in practice. As shown in Figure 1, we notice a
33 severe false-alarm problem experienced by previous methods when the semantic consensus does not
34 exist, which may prevent those methods from being useful in various applications that cannot provide
35 accurate vision-language pairs. We argue that the current R-VOS task is not completely defined with
36 the assumption that the referred object always exists in the video.

37 Even when semantic consensus exists in the given video-language pairs, it is still challenging to locate
38 the correct object in the video due to the multimodal nature of the R-VOS task. Recently, MTTR [1]

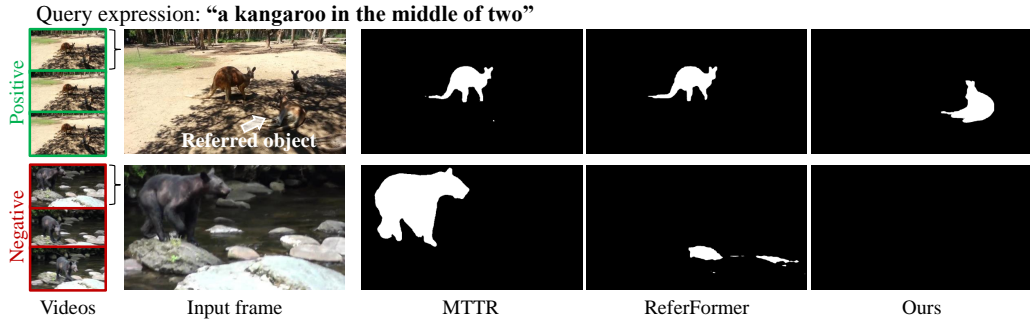


Figure 1: Illustration of the new R^2 -VOS task. A linguistic expression is given to query a set of videos without the semantic consensus assumption. Videos containing the referred object by the expression are **positive**, otherwise **negative**. Unlike the previous R-VOS setting that assumes all target videos are positive to the query expression, the new R^2 -VOS task is required to discriminate positive and negative text-video pairs, and further segment object masks for all frames in positive videos or treat entire negative videos as backgrounds. Compared to the previous state-of-the-art R-VOS methods, MTTR [1] and ReferFormer [44], our method not only discriminates negative videos better but also shows a superior disambiguating ability between visually similar objects in positive videos.

39 employs a multimodal transformer encoder to learn a joint representation of the linguistic expression
 40 and video, and then obtains the referred object by ranking all objects in the video. ReferFormer
 41 [44] follows the image-level method, ReTR [19], to adopt the linguistic expression as a query to
 42 the transformer decoder to avoid redundant ranking of all objects. However, these latest methods
 43 suffer from semantic misalignment of the segmented object and the linguistic expression, even with
 44 sophisticated components employed. As shown in Figure 1, the segmented objects by MTTR and
 45 ReferFormer are often not the object referred to by the linguistic expression.

46 In this paper, we seek to investigate the semantic alignment problem between visual and linguistic
 47 modalities in referring video segmentation. We extend the current task definition of R-VOS [13] to
 48 accept both paired and unpaired video and language inputs. This new task, which we term Robust
 49 R-VOS (R^2 -VOS), overcomes the current limitation of the R-VOS task by additionally considering
 50 the semantic alignment of input video to referring expression. We reveal that this task is essentially
 51 related to two problems that are interrelated [31]: the R-VOS problem as the **primary** problem of
 52 segmenting mask sequences from videos with referring texts, and its **dual** problem of reconstructing
 53 text expressions from videos with object masks. By linking the primary and dual problems, we
 54 introduce a text-video-text cycle and a corresponding relational consistency constraint, which can
 55 enforce the semantic consensus between the text query and segmented mask to improve the primary
 56 task. In practice, naively conducting cyclic training of the text-video-text cycle will lead to a two-
 57 stage regime and significantly increasing costs. We address this problem by incorporating an early
 58 grounding scheme, serving as a **proxy**, to efficiently model the two tasks in a parallel manner. In
 59 addition, we discriminate the semantic misalignment between the video and text by assessing the cycle
 60 consistency between the original and reconstructed texts, thus alleviating the false-alarm problem.
 61 Our contributions can be summarized as:

- 62 • We notice a severe false-alarm problem faced by previous R-VOS methods with unpaired
 63 inputs. To investigate the robustness of current referring segmentation models, we introduce
 64 the R^2 -VOS task that accepts unpaired video and text as inputs.
- 65 • We propose a pipeline that jointly optimizes the primary referring segmentation and dual
 66 expression reconstruction task and introduces a relational cycle consistency constraint to
 67 enhance the semantic alignment between visual and textual modalities.
- 68 • Our method surpasses previous state-of-the-art methods on Ref-Youtube-VOS, Ref-DAVIS,
 69 and R^2 -Youtube-VOS dataset in terms of both performance and speed.

70 2 Related Works

71 **Vision and language representation learning.** There have been a long line of studies on how to
 72 learn better vision-language representation, e.g., multimodal attention [30, 50, 8, 3], fusion scheme
 73 [7, 14, 15, 51], multi-step reasoning [47, 10] and pretraining [37, 5, 17]. KAC Net [2] leverages

74 knowledge-aided consistency constraints to enhance semantic alignment for weakly supervised phrase
 75 grounding. A structure-preserving constraint [42] is proposed to preserve some intra-modal properties
 76 when learning vision-language representation for image-text retrieval.

77 **Referring video object segmentation.** R-VOS is a novel task that aims to segment an object across
 78 frames given a linguistic description. URVOS [39] is the first unified R-VOS framework with a
 79 cross-modal attention and a memory attention module, which largely improves R-VOS performance.
 80 ClawCraneNet [21] leverages cross-modal attention to bridge the semantic correlation between textual
 81 and visual modalities. ReferFormer [44] and MTTR [1] are two latest works that utilize transformers
 82 to decode or fuse multimodal features. ReferFormer [44] employs a linguistic prior to the transformer
 83 decoder to focus on the referred object. MTTR [1] leverages a multimodal transformer encoder
 84 to fuse linguistic and visual features. Different from other vision-language tasks, e.g., image-text
 85 retrieval [25, 26, 32] and video question answering [18, 40], R-VOS needs to construct object-level
 86 multimodal semantic consensus in a dense visual representation.

87 3 R²-VOS

88 3.1 Task Definition

89 We introduce a novel task, robust referring video segmentation (R²-VOS), which aims to predict
 90 mask sequences $\{M_o\}$ for an unconstrained video set $\{V\}$ given a language expression E_o of an
 91 object o . Different from the previous R-VOS setup, the queried video V is not required to contain
 92 the referred object by expression E_o . We define a video V and an expression E_o to have **semantic**
 93 **consensus** if the object o appears in V , and the video is **positive** with respect to E_o , otherwise it is
 94 **negative**. The R²-VOS task is extended to discriminate positive and negative videos, and predict
 95 masks M_o of object o for positive videos and treat all frames in the negative videos as background.

96 3.2 Problem Analysis

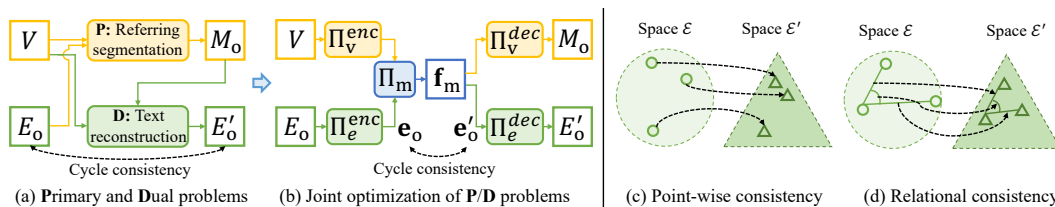


Figure 2: Problem analysis. (a) R²-VOS introduces the **Primary** problem of referring segmentation and the **Dual** problem of text reconstruction for positive videos. The **P/D** problems are connected in a cycle path from original expression E_o to reconstructed expression E'_o . (b) The cycle consistency between the original and reconstructed embeddings (e_o and e'_o) can benefit to optimize the **P** problem. We enable the joint optimization for cycle consistency with a cross-modal **proxy** f_m defined between all single-modal operations (i.e., Π_v^{enc} , Π_e^{enc} , Π_v^{dec} and Π_e^{dec}). (c) Point-wise consistency is not suitable in R²-VOS because the mapping between \mathcal{E} and \mathcal{E}' are not necessarily bijective. **An object can be referred by various textual expressions.** (d) Instead, we apply a relational consistency to preserve distances and angles.

97 **Primary and dual problems for R²-VOS.** The referring segmentation can be formulated as the
 98 maximum *a posteriori* estimation problem of $p(M_o|V, E_o)$. By applying the Bayes rule, we obtain:

$$p(M_o|V, E_o) \sim p(E_o|V, M_o)p(M_o|V) \quad (1)$$

99 As the prior $p(M_o|V)$ is not affected by the expression E_o , we consider maximizing $p(E_o|V, M_o)$
 100 as a dual problem of the referring segmentation (primary problem), which is to reconstruct the text
 101 expression given the video and object masks. We note that for negative videos, $p(E_o|V, M_o)$
 102 is undefined because the mask M_o is empty. Thus, we only investigate the dual problem for positive
 103 videos. The primary problem and the dual problem can be connected in a cycle path, i.e., from the
 104 original expression E_o to the reconstructed expression E'_o through positive video queries, as shown
 105 in Figure 2 (a). We believe that the cycle constraint benefits to optimize the primary problem by
 106 enhancing the semantic consensus.

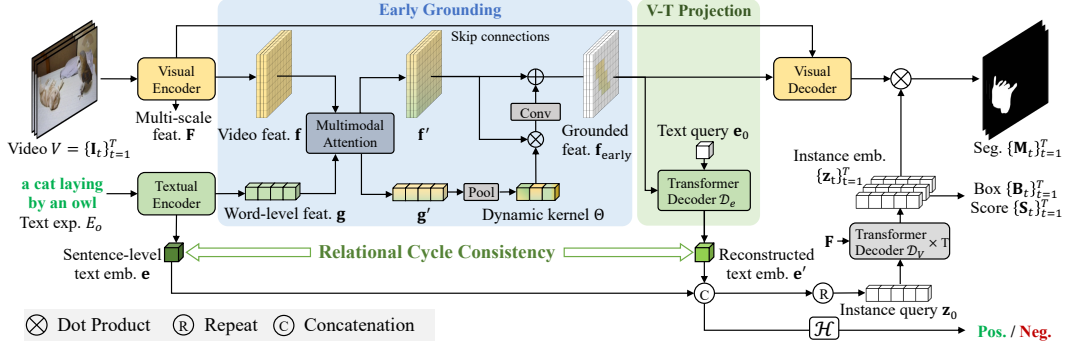


Figure 3: Overview of the proposed model. Given a video clip $V = \{\mathbf{I}_t\}_{t=1}^T$ and a textual expression E_o referring object o , we first extract video feature and text feature separately, then fuse them in the early grounding module to obtain the visual representation $\mathbf{f}_{\text{early}}$ of the referred object o . Then we project $\mathbf{f}_{\text{early}}$ to a textual space to be \mathbf{e}' and add the relational cycle constraint with the original text embedding \mathbf{e} . The final segmentation is obtained by dynamic convolutions with video features from the visual decoder and dynamic weights from the fused text embeddings. The semantic consensus of input pairs is discriminated to be positive or negative by assessing the consistency between \mathbf{e} and \mathbf{e}' .

107 In practice, we study the cycle consistency between the original textual embedding space \mathcal{E} and the
 108 transformed **textual** embedding space \mathcal{E}' induced by positive videos. By definition, the path from the
 109 original text embedding \mathbf{e}_o to the reconstructed text embedding \mathbf{e}'_o is modulated with **cross-modal**
 110 interactions between video and text. Thus, to link the primary and dual problem and enable the joint
 111 optimization, we introduce a cross-modal **intermediate feature** \mathbf{f}_m to convey information of both the
 112 input of the primary problem (V, E_o) and the dual problem (V, M_o) , as shown in Figure 2 (b). \mathbf{f}_m
 113 is defined between the encoder and decoder stages of single-modal operations, i.e., $\Pi_v^{\text{enc}}, \Pi_e^{\text{enc}}, \Pi_v^{\text{dec}}$,
 114 Π_e^{dec} , to only focus on the multi-modal interaction.

115 **Relational cycle consistency.** A key observation for cycle consistency between \mathcal{E} and \mathcal{E}' is that the
 116 mapping between them is not necessarily bijective, as there could be multiple textual descriptions for
 117 the same object. Thus, naively adding point-wise consistency, i.e., $\mathbf{e}_o = \mathbf{e}'_o, \forall \mathbf{e}_o \in \mathcal{E}$ will collapse
 118 the feature space to a sub-optimal solution. Instead, we take inspiration from relational knowledge
 119 distillation [33], and introduce relational cycle consistency for \mathcal{E} and \mathcal{E}' . The relational cycle
 120 consistency is to preserve the structure of the feature space rather than exact point-wise consistency,
 121 as illustrated in Figure 2 (c) and (d). Mathematically, the structure-preserving property is defined as
 122 isometric and conformal constraints to preserve pair-wise distance and angles for $\mathbf{e} \in \mathcal{E}$ and $\mathbf{e}' \in \mathcal{E}'$:

$$|\mathbf{e}_1 - \mathbf{e}_2| = |\mathbf{e}'_1 - \mathbf{e}'_2| \quad (2)$$

$$\angle(\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3) = \angle(\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3), \quad (3)$$

123 where $|\cdot|$ and $\angle(\cdot)$ denote distance and angle metrics.

124 4 Method

125 In this section, we elaborate our $\text{R}^2\text{-VOS}$ framework with the relational consistency, which mainly
 126 consists of four parts: feature extraction, early grounding as a medium, video-text (V-T) projection for
 127 text reconstruction, and mask decoding for final segmentation, as shown in Figure 3. We first extract
 128 the video feature \mathbf{f} , word-level text feature \mathbf{g} , and sentence-level text embedding \mathbf{e} . On the one hand,
 129 to model the primary segmentation problem of maximizing $p(M_o|V, E_o)$, we enable the multimodal
 130 interaction in the early grounding module to generate the grounded feature $\mathbf{f}_{\text{early}}$. $\mathbf{f}_{\text{early}}$ coarsely
 131 locates the referred object o and filters out irrelevant features, which serves as a medium linking
 132 the primary segmentation and dual text reconstruction problem. The final mask M_o is obtained by
 133 dynamic convolution [4] on the decoded visual feature maps, with kernels learned from instance
 134 embedding $\{\mathbf{z}_t\}_{t=1}^T$. On the other hand, to model the dual text reconstruction problem of maximizing
 135 $p(E_o|V, M_o)$, we utilize the grounded video feature $\mathbf{f}_{\text{early}}$ as the alternative of V and M_o , since
 136 $\mathbf{f}_{\text{early}}$ conveys contextual video clues of object o . In this way, we enable the parallel optimization of
 137 the primary and dual problem by relating them to $\mathbf{f}_{\text{early}}$. Specifically, we employ a V-T projection
 138 module to project $\mathbf{f}_{\text{early}}$ onto a reconstructed text embedding \mathbf{e}' . We add relational constraint between

139 e' and e to enforce the semantic alignment between the segmented mask and expression for positive
 140 videos. In addition, we introduce a semantic consensus discrimination head $\mathcal{H}(e, e')$ to assess the
 141 consistency between original and reconstructed text embeddings, discriminating the alignment of
 142 multimodal semantics and identifying negative videos.

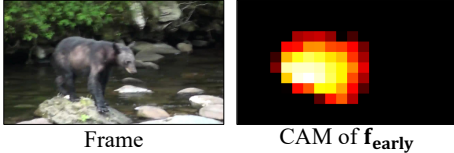
143 4.1 Single-modal Feature Extraction

144 **Visual encoder.** Following previous methods [1, 44, 43], we build the visual encoder with a visual
 145 backbone and a deformable transformer encoder [52] on top of it. The extracted features from the
 146 backbone are flattened, projected to a lower dimension, added with positional encoding [12], and
 147 then fed into a deformable transformer encoder [52] similar to the previous method [44]. We denote
 148 the multi-scale output of the transformer encoder as \mathbf{F} and the low-resolution visual feature map from
 149 the backbone as \mathbf{f} , where $\mathbf{f} \in \mathbb{R}^{T \times C_v \times \frac{H}{32} \times \frac{W}{32}}$, C_v is the feature channel, T is the video length and H
 150 and W are the original image size.

151 **Textual encoder.** We leverage a pre-trained linguistic model RoBERTa [27] to map the input textual
 152 expression E_o to a textual embedding space. The textual encoder extracts a sequence of word-level
 153 text feature $\mathbf{g} \in \mathbb{R}^{C_e \times L}$ and a sentence-level text embedding $e \in \mathbb{R}^{C_e \times 1}$, where C_e and L are the
 154 dimension of linguistic embedding space and the expression length respectively.

155 4.2 Early Grounding

156 **a black bear standing on a rock in a stream**



157
158
159
160
161
162 Figure 4: Visualization of channel activation map (CAM) of $\mathbf{f}_{\text{early}}$.

163
164
165 early stage. As shown in the blue part of Figure 3, we first enable the multimodal interaction between
 166 video and text features, then apply the dynamic convolution with kernels learned from text feature
 167 to discriminate the object-level semantics. In particular, multi-head cross-attention (MCA) [41] is
 168 leveraged to conduct the multimodal interaction:

$$169 \quad \mathbf{h}_f = \text{LN}(\text{MCA}(\mathbf{f}, \mathbf{g}) + \mathbf{f}) \quad \mathbf{f}' = \text{LN}(\text{FFN}(\mathbf{h}_f) + \mathbf{h}_f) \quad (4)$$

$$\mathbf{h}_g = \text{LN}(\text{MCA}(\mathbf{g}, \mathbf{f}) + \mathbf{g}) \quad \mathbf{g}' = \text{LN}(\text{FFN}(\mathbf{h}_g) + \mathbf{h}_g), \quad (5)$$

170 where $\text{MCA}(\mathbf{X}, \mathbf{Y}) = \text{Attention}(\mathbf{W}^Q \mathbf{X}, \mathbf{W}^K \mathbf{Y}, \mathbf{W}^V \mathbf{Y})$. \mathbf{W} represents learnable weight. LN
 171 and FFN denote layer normalization and feed-forward network respectively. The text feature \mathbf{g}' is
 172 further pooled to a fixed length, and followed by a fully-connected layer to form the dynamic kernels
 173 $\Theta = \{\theta_i\}_{i=1}^K$. K is the kernel number and $\theta_i \in \mathbb{R}^{C \times 1}$. The dynamic kernels are applied separately
 174 to video feature $\mathbf{f}' \in \mathbb{R}^{C \times THW}$ to form the $\mathbf{f}_{\text{early}} \in \mathbb{R}^{C \times THW}$

$$\mathbf{f}_{\text{early}} = \text{BN}(\varphi(\theta_1^T \mathbf{f}' \oplus \dots \oplus \theta_K^T \mathbf{f}') + \mathbf{f}'), \quad (6)$$

175 where \oplus is the concatenation in channel dimension and $\varphi(\cdot)$ is a convolution to reduce the feature
 176 dimension. BN denotes batch normalization.

177 4.3 Text Reconstruction

178 **V-T projection.** We leverage a transformer decoder \mathcal{D}_E as textual decoder to transform the visual
 179 representation of the referred object into the textual space. As shown in Figure 3, a learnable text
 180 query $\mathbf{e}_0 \in \mathbb{R}^{C_e \times 1}$ is employed to query the $\mathbf{f}_{\text{early}}$. The output of the transformer decoder is the
 181 reconstructed text embedding $e' = \mathcal{D}_E(\mathbf{f}_{\text{early}}, \mathbf{e}_0) \in \mathbb{R}^{C_e \times 1}$.

182 4.4 Referring Segmentation

183 **Mask segmentation.** Similar to previous methods [44, 1, 11], we leverage deformable transformer
 184 decoders with dynamic convolution to segment the object masks. As shown in Figure 3, we first fuse

185 the reconstructed text embedding e' to text embedding e . The fused text embedding e is then repeated
 186 N times to form the instance query [43] $\mathbf{z}_0 \in \mathbb{R}^{C_q \times N}$, where C_q is the dimension of instance query
 187 and N is the instance query number. We then use $T \times$ deformable transformer decoders \mathcal{D}_V with
 188 shared weights to decode the instance embeddings $\mathbf{z}_t \in \mathbb{R}^{C_q \times N}$ for each frame, i.e., $\mathbf{z}_t = \mathcal{D}_V(\mathbf{F}_t, \mathbf{z}_0)$.
 189 \mathbf{F}_t is the multiscale visual feature from visual encoder at time t . A dynamic kernel \mathbf{w}_t is further
 190 learned from the instance embedding \mathbf{z}_t . The final feature map $\mathbf{f}_{\text{out},t} \in \mathbb{R}^{C \times H \times W}$ is obtained by
 191 fusing low-level features from the feature pyramid network [23] in the visual decoder. The mask
 192 prediction $\mathbf{M}_t \in \mathbb{R}^{N \times H \times W}$ can be computed by $\mathbf{M}_t = \mathbf{w}_t^T \mathbf{f}_{\text{out},t}$.

193 **Auxiliary heads.** We build a set of auxiliary heads to obtain the final object masks across frames. In
 194 particular, a box head, a scoring head and a semantic consensus discrimination head are leveraged to
 195 predict the bounding boxes $\mathbf{B}_t \in \mathbb{R}^{N \times 4}$, confidence scores $\mathbf{S}_t \in \mathbb{R}^{N \times 1}$ and the alignment degree of
 196 multimodal semantics $A \in \mathbb{R}$. The box and scoring head are two fully-connected layers upon the
 197 instance embedding \mathbf{e}_t . The semantic consensus discrimination head $\mathcal{H}(e, e')$ consists of two fully-
 198 connected layers upon the text embeddings $e \oplus e'$. Note that A represents the semantic alignment in
 199 the entire video rather a single frame, since the expression is a video-level description.

200 4.5 Loss Function

201 The loss function of our method can be boiled down to three parts:

$$\mathcal{L} = \lambda_{\text{text}} \mathcal{L}_{\text{text}} + \lambda_{\text{segm}} \mathcal{L}_{\text{segm}} + \lambda_{\text{align}} \mathcal{L}_{\text{align}}, \quad (7)$$

202 where $\mathcal{L}_{\text{text}}$, $\mathcal{L}_{\text{segm}}$, and $\mathcal{L}_{\text{align}}$ are losses for text reconstruction, referring segmentation and semantic
 203 consensus discrimination respectively. A ground-truth semantic alignment $\hat{A} = \{0, 1\}$ is utilized
 204 to discriminate positive and negative pairs. The $\mathcal{L}_{\text{align}}$ is simply a cross-entropy loss between the
 205 predicted alignment A and ground-truth \hat{A} . The other two terms are computed as follows:

206 **Loss for text reconstruction.** Given the text embedding e and reconstructed text embedding e' , we
 207 employ a relational constraint to impose the cycle consistency between e and e' . We calculate the
 208 loss by

$$\mathcal{L}_{\text{text}} = \mathbb{1}(\hat{A}) \cdot (\mathcal{L}_{\text{dist}} + \lambda_{\text{angle}} \mathcal{L}_{\text{angle}}), \quad (8)$$

209 where the indicator function $\mathbb{1}(\hat{A}) = 1$ if the alignment indicates the referred object exists in the
 210 video, otherwise 0, λ_{angle} is a hyperparameter balancing the distance loss $\mathcal{L}_{\text{dist}}$ and angle loss
 211 $\mathcal{L}_{\text{angle}}$. We elaborate these two losses according to the relational cycle consistency Equation 2.
 212 Let $\mathcal{X}^n = \{(x_1, \dots, x_n) | x_i \in \mathcal{X}\}$ denote a set of n -tuples, $\Phi^n = \{(\mathbf{x}, \mathbf{x}') | \mathbf{x} \in \mathcal{X}^n, \mathbf{x}' \in \mathcal{X}'^n\}$
 213 denote a set of pairs consisting of two n -tuples of distinct elements from two different sets \mathcal{X} and
 214 \mathcal{X}' . Specifically, the distance-based and angle-based relations relate text embeddings of 2-tuple and
 215 3-tuple respectively, i.e., $\Phi^2 = \{(\mathbf{x}, \mathbf{x}') | \mathbf{x} = (\mathbf{e}_i, \mathbf{e}_j), \mathbf{x}' = (\mathbf{e}'_i, \mathbf{e}'_j), i \neq j\}$ and $\Phi^3 = \{(\mathbf{x}, \mathbf{x}') | \mathbf{x} =$
 216 $(\mathbf{e}_i, \mathbf{e}_j, \mathbf{e}_k), \mathbf{x}' = (\mathbf{e}'_i, \mathbf{e}'_j, \mathbf{e}'_k), i \neq j \neq k\}$. Then the losses are given by:

$$\mathcal{L}_{\text{dist}} = \sum_{(\mathbf{x}, \mathbf{x}') \in \Phi^2} l_\delta(\phi_D(\mathbf{x}), \phi_D(\mathbf{x}')), \quad \phi_D(\mathbf{x}) = \frac{1}{\mu(\mathbf{x})} \|\mathbf{e}_i - \mathbf{e}_j\|_2, \quad (9)$$

$$\mathcal{L}_{\text{angle}} = \sum_{(\mathbf{x}, \mathbf{x}') \in \Phi^3} l_\delta(\phi_\angle(\mathbf{x}), \phi_\angle(\mathbf{x}')), \quad \phi_\angle(\mathbf{x}) = \cos \angle(\mathbf{e}_i, \mathbf{e}_j, \mathbf{e}_k), \quad (10)$$

217 where $\mu(\mathbf{x}) = \sum_{\mathbf{x}=(x_1, x_2) \in \mathcal{X}^2} \frac{\|x_1 - x_2\|_2}{|\mathcal{X}^2|}$ is the average distance function, and the Huber loss
 218 $l_\delta(x, x') = \frac{1}{2}(x - x')^2$ if $|x - x'| \leq 1$, otherwise $|x - x'| - \frac{1}{2}$.

219 **Loss for referring segmentation.** Given a set of predictions $\mathbf{y} = \{\mathbf{y}_i\}_{i=1}^N$ and ground-truth $\hat{\mathbf{y}}$,
 220 where $\mathbf{y}_i = \{\mathbf{B}_{i,t}, \mathbf{S}_{i,t}, \mathbf{M}_{i,t}\}_{t=1}^T$ and $\hat{\mathbf{y}} = \{\hat{\mathbf{B}}_t, \hat{\mathbf{S}}_t, \hat{\mathbf{M}}_t\}_{t=1}^T$, we search for an assignment $\sigma \in \mathcal{P}_N$
 221 with the highest similarity where \mathcal{P}_N is a set of permutations of N elements ($\hat{\mathbf{y}}$ is padded with \emptyset).
 222 The similarity can be computed as

$$\mathcal{L}_{\text{match}}(\mathbf{y}_i, \hat{\mathbf{y}}) = \lambda_{\text{box}} \mathcal{L}_{\text{box}} + \lambda_{\text{conf}} \mathcal{L}_{\text{conf}} + \lambda_{\text{mask}} \mathcal{L}_{\text{mask}}, \quad (11)$$

223 where λ_{box} , λ_{conf} , and λ_{mask} are weights to balance losses. Following previous works [6, 43], we
 224 leverage a combination of Dice [20] and BCE loss as $\mathcal{L}_{\text{mask}}$, focal loss [24] as $\mathcal{L}_{\text{conf}}$, and GIoU
 225 [38] and L1 loss as \mathcal{L}_{box} . The best assignment $\hat{\sigma}$ is solved by Hungarian algorithm [16]. Given
 226 the best assignment $\hat{\sigma}$, the segmentation loss between ground-truth and predictions is defined as
 227 $\mathcal{L}_{\text{segm}} = \mathbb{1}(\hat{A}) \cdot \mathcal{L}_{\text{match}}(\mathbf{y}, \hat{\mathbf{y}}_{\hat{\sigma}(i)})$.

228 4.6 Inference

229 During inference, we select the candidate with the highest confidence to predict the final masks:

$$\{\bar{\mathbf{M}}_t\}_{t=1}^T = \{\mathbb{1}(A) \cdot \mathbf{M}_{\bar{s},t}\}_{t=1}^T, \quad \bar{s} = \underset{i}{\operatorname{argmax}}\{\mathbf{S}_{i,1} + \dots + \mathbf{S}_{i,T}\}_{i=1}^N, \quad (12)$$

230 where $\{\bar{\mathbf{M}}_t\}_{t=1}^T$ is the masks of referred object. $\mathbf{S}_{i,t}$ and $\mathbf{M}_{i,t}$ represent the i -th candidate in \mathbf{S}_t and
231 \mathbf{M}_t respectively. \bar{s} is the slot with the highest confidence to be the target object. We use $\mathbb{1}(A)$ to filter
232 out predictions in negative videos to mitigate false alarm. $\mathbb{1}(A) = 1$ if $A > 0.5$, else 0.

233 5 Experiment

234 5.1 Dataset and Metrics

235 **Dataset.** We conduct experiments on three datasets: Ref-Youtube-VOS, Ref-DAVIS and R²-
236 Youtube-VOS. Ref-Youtube-VOS [39] is a large-scale benchmark that has 3,978 videos with about
237 15k language descriptions. There are 3,471 videos with 12,913 expressions in the training set and 507
238 videos with 2,096 expressions in the validation set. Ref-DAVIS-17 [13] contains 90 videos with 1,544
239 expressions, including 60 and 30 videos for training and validation respectively. R²-Youtube-VOS
240 is our newly proposed evaluation dataset: it extends the Ref-Youtube-VOS validation set with each
241 linguistic expression to query a positive video (the same one as Ref-Youtube-VOS) and a negative
242 video. To make each video can be picked as a negative video, we randomly shuffle the original video
243 set and constrain all negative text-video pair unrelated.

244 **Metrics.** We employ commonly-used region similarity \mathcal{J} and contour accuracy \mathcal{F} [36] for con-
245 ventional Ref-Youtube-VOS and Ref-DAVIS-17 benchmarks. For the proposed R²-Youtube-VOS
246 task, we additionally introduce a new metric $\mathcal{R} = 1 - \frac{\sum_{M \in \mathcal{M}_{neg}} |M|}{\sum_{M \in \mathcal{M}_{pos}} |M|}$ to evaluate the degree of object
247 false alarm in negative videos, where \mathcal{M}_{neg} and \mathcal{M}_{pos} are the sets containing segmented masks in
248 negative and positive videos respectively. $|M|$ denotes the foreground area of mask M . The total
249 foreground area of positive videos $\sum_{M \in \mathcal{M}_{pos}} |M|$ serves as a normalization term. Ideally, a model
250 should treat all the negative videos as backgrounds with $\mathcal{R} = 1$.

251 5.2 Implementation Details

252 Following previous methods [6, 44], our model is first pre-trained on Ref-COCO+/g dataset [49, 31]
253 and then finetuned on Ref-Youtube-VOS. The model is trained for 6 epochs with a learning rate
254 multiplier of 0.1 at the 3rd and the 5th epoch. The initial learning rate is 1e-4 and a learning rate
255 multiplier of 0.5 is applied to the backbone. We adopt a batchsize of 8 and an AdamW [29] optimizer
256 with weight decay 1×10^{-4} . Following convention [44, 1], the evaluation on Ref-DAVIS directly
257 uses models trained on Ref-Youtube-VOS without re-training. All images are cropped to have the
258 longest side of 640 pixels and the shortest side of 360 pixels during evaluation. [The window size is](#)
259 [set to 5 for all backbones. We create negative pairs by shuffling positive pairs in each batch.](#) Our
260 method is implemented with PyTorch [34].

261 5.3 Main Results

262 We compare our method with state-of-the-art R-VOS methods on Ref-Youtube-VOS and Ref-DAVIS-
263 17 in Table 1, and R²-VOS task in Table 2.

264 **Comparison on Ref-Youtube-VOS.** In Table 1, we first compare our method on Ref-Youtube-VOS.
265 For results of ResNet [9] backbone, our method achieves 57.3 $\mathcal{J}\&\mathcal{F}$ which outperforms the latest
266 method ReferFormer [44] by 1.7 $\mathcal{J}\&\mathcal{F}$. In addition, our method runs at 30 FPS compared to 22 FPS
267 of state-of-the-art ReferFormer (FPS is measured using single NVIDIA P40 with *batchsize* = 1).
268 For results of Swin-Transformer [28, 28] backbones, our method achieves 60.2 $\mathcal{J}\&\mathcal{F}$ and 61.3 $\mathcal{J}\&\mathcal{F}$
269 with Swin-Tiny and Video-Swin-Tiny backbones respectively, which outperforms ReferFormer [44]
270 and MTTR [1] by a clear margin. [More analysis is available in the additional appendix A.1.](#)

271 **Comparison on Ref-DAVIS-17.** Our method achieves 59.7 $\mathcal{J}\&\mathcal{F}$ on Ref-DAVIS-17 dataset, which
272 outperforms ReferFormer by 1.2 $\mathcal{J}\&\mathcal{F}$.

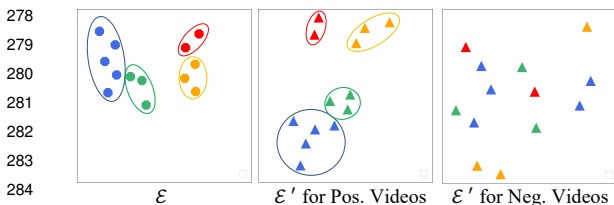
Method	Backbone	Ref-Youtube-VOS			Ref-DAVIS-17		
		$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}
Spatial Visual Backbone							
CMSA [48]	ResNet-50	34.9	33.3	36.5	34.7	32.2	37.2
CMSA + RNN [48]	ResNet-50	36.4	34.8	38.1	40.2	36.9	43.5
URVOS [39]	ResNet-50	47.2	45.3	49.2	51.5	47.3	56.0
PMINet [6]	ResNet-101	53.0	51.5	54.5	-	-	-
CITD [22]	ResNet-101	56.4	54.8	58.1	-	-	-
ReferFormer* [44]	ResNet-50	55.6	54.8	56.5	58.5	55.8	61.3
Ours	ResNet-50	57.3	56.1	58.4	59.7	57.2	62.4
ReferFormer* [44]	Swin-T	58.7	57.6	59.9	-	-	-
Ours	Swin-T	60.2	58.9	61.5	-	-	-
Spatio-temporal Visual Backbone							
MTTR* [1]	Video-Swin-T	55.3	54.0	56.6	-	-	-
ReferFormer* [44]	Video-Swin-T	59.4	58.0	60.9	-	-	-
Ours	Video-Swin-T	61.3	59.6	63.1	-	-	-

Table 1: Comparison to state-of-the-art R-VOS methods on Ref-Youtube-VOS and Ref-DAVIS-17 val set. * indicates results imported from preprints.

Method	Backbone	$\mathcal{J}\&\mathcal{F}$ & \mathcal{R}	\mathcal{J}	\mathcal{F}	\mathcal{R}
ReferFormer* [44]	ResNet-50	47.3	54.8	56.5	30.6
Ours	ResNet-50	69.5	56.1	58.4	94.1
MTTR* [1]	Video-Swin-T	40.0	55.9	58.1	5.9
ReferFormer* [44]	Video-Swin-T	49.1	58.0	60.9	28.5
Ours	Video-Swin-T	72.8	59.6	63.1	95.7

Table 2: Comparison to state-of-the-art R-VOS methods on R²-Youtube-VOS.

273 **Comparison on R²-VOS.** As shown in Table 2, the state-of-the-art R-VOS methods, ReferFormer
274 and MTTR suffer from a low \mathcal{R} metric which measures the false-alarm problem when the semantic
275 consensus of the input text-video pair does not hold. Compared to the severe false alarm of previous
276 R-VOS methods, our model successfully mitigates the false alarm of the model, thanks to the proposed
277 multimodal cycle consistency constraint and semantic consensus discrimination head.

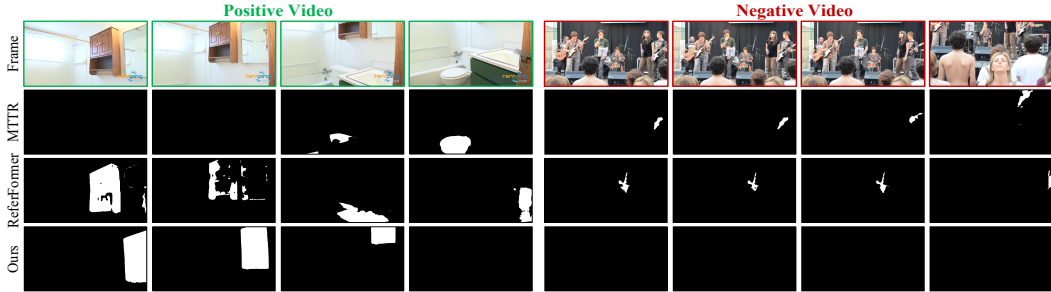


280 Figure 5: Visualization of text embedding spaces. Dots
281 represent original text embeddings in \mathcal{E} , and triangles
282 represent reconstructed ones in \mathcal{E}' induced by positive
283 and negative videos respectively. Elements in the same
284 color belong to the same object. Note that an object
285 can have multiple text descriptions. The structure of \mathcal{E}'
286 is well preserved from \mathcal{E} for positive videos (ellipses
287 bound embeddings of same objects), while it is not
288 preserved for negative videos.

289 reconstructed text embedding spaces for both positive and negative videos. As shown in Figure 5, we
290 notice that, for embeddings of positive videos, they preserve relative relations well, while for negative
291 videos, the reconstructed embeddings have a random pattern in the space.

298 5.4 Ablation Study

299 **Module effectiveness.** To investigate the effectiveness of different components in our method,
300 we conduct experiments with the ResNet-50 backbone on R²-Youtube-VOS dataset. We build a
301 transformer-based baseline model and equip our proposed components step-by-step. As shown in
302 Table 3, the baseline model achieves 52.4 $\mathcal{J}\&\mathcal{F}$. Then, we add our proposed components step-by-
303 step to demonstrate the module effectiveness. After employing the early grounding module, the
304 performance boosts to 55.5 $\mathcal{J}\&\mathcal{F}$ and the cycle-consistency constraint brings another 1.4 $\mathcal{J}\&\mathcal{F}$ gain.
305 Since the reconstructed text embedding is generated with visual features injected, we consider it can



Expression: **the mirror in the bathroom is to the right of the wood cabinet**

Figure 6: Qualitative comparison to the state-of-the-art R-VOS method on the R²-VOS task.

Components	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	\mathcal{R}
Baseline	52.4	51.9	52.8	34.9
+EG	55.5 ^{+3.1}	54.4	56.5	32.9 ^{-2.0}
+CC	56.9 ^{+4.5}	55.7	58.1	94.0 ^{+59.1}
+FT	57.3^{+4.9}	56.1	58.4	94.1^{+59.2}

Table 3: **Impact of different components in our method.** EG: Early grounding, CC: Consistency constraint, FT: Fusing text embeddings.

Query Number	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	\mathcal{R}
1	54.9	54.2	55.6	94.7
5	57.3	56.1	58.4	94.1
9	57.0	56.8	57.2	93.5

Table 5: **Impact of the query number.**

Constraint	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	\mathcal{R}
None	55.5	54.4	56.5	32.9
PW	54.4 ^{-1.1}	53.3	55.5	88.7 ^{+55.8}
RA	56.7 ^{+1.2}	55.5	57.9	93.6 ^{+60.7}
RD	56.4 ^{+0.9}	55.2	57.6	90.4 ^{+57.5}
RD+RA	56.9^{+1.4}	55.7	58.1	94.0^{+61.1}

Table 4: **Impact of the cycle consistency constraint.** PW: Point-wise. RA: Relational angle. RD: Relational distance.

Window Size	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	\mathcal{R}
1	53.5	53.0	54.0	89.2
3	56.8	56.5	57.1	92.1
5	57.3	56.1	58.4	94.1

Table 6: **Impact of the window size.**

306 encode some visual information, thus augmenting the original text embedding. By using the fused
 307 text embedding as instance query, we achieve our best performance of 57.3 $\mathcal{J}\&\mathcal{F}$.

308 **Consistency constraint.** We conduct experiments to ablate the influence of cycle-consistency
 309 constraints. As shown in Table 4, utilizing point-wise consistency constraint will lead to a performance
 310 drop compared to the setting without cycle constraint. We consider the point-wise constraint may
 311 force an injective mapping from the textual domain to the visual domain. However, the mapping can
 312 be a many-to-one function for R-VOS, i.e., each object corresponds to multiple textual descriptions.
 313 In addition, since the early grounding leverages the text feature to locate the referred object, if we use
 314 the direct point-wise constraint to form reconstructed text embedding, it will encourage the network
 315 to memorize the text feature in the $\mathbf{f}_{\text{early}}$ and result in a collapse for text reconstruction. Table 4
 316 shows that sole relational angle constraint can bring 1.2 $\mathcal{J}\&\mathcal{F}$ gain, and it can be slightly improved
 317 with 1.4 $\mathcal{J}\&\mathcal{F}$ gain by jointly using relational angle and distance constraint.

318 **Instance query number.** Although only one referral is involved for each frame in R-VOS task,
 319 to help the network optimization, we employ more than one instance query to each video. Table 5
 320 indicates that a query number of 5 brings the best result.

321 **Frame number.** Since R-VOS gives a text that describes an object over a period of time, temporal
 322 information is vital to segment accurate and temporally-consistent results. We ablate on the best
 323 window size of input videos during training. As shown in Table 6, we notice that the performance
 324 improves as the window size increases and a window size of 5 brings the best result of 57.3 $\mathcal{J}\&\mathcal{F}$.

325 6 Conclusion

326 In this paper, we investigate the semantic misalignment problem in R-VOS task. A pipeline jointly
 327 models the referring segmentation and text reconstruction problem, equipped with a relational cycle
 328 consistency constraint, is introduced to discriminate and enhance the semantic consensus between
 329 visual and textual modalities. To evaluate the model robustness, we extend the R-VOS task to
 330 accept unpaired inputs and collect a corresponding R²-Youtube-VOS dataset. We observe a severe
 331 false-alarm problem suffered from previous methods on R²-Youtube-VOS while ours accurately
 332 discriminates unpaired inputs and segments high-quality masks for paired inputs. Our method
 333 achieves state-of-the-art performance on Ref-DAVIS17, Ref-Youtube-VOS, and R²-VOS dataset. We
 334 believe that, with unpaired inputs, R²-VOS is a more general setting of referring video segmentation,
 335 which can shed light on a new direction to investigate the robustness of referring segmentation.

336 A Additional Appendix

337 A.1 More Quantitative Result Analysis

338 Under the same ResNet-50 backbone, our method achieves 57.3 $\mathcal{J}\&\mathcal{F}$, 94.1 \mathcal{R} and 30 FPS compared
339 to the 55.6 $\mathcal{J}\&\mathcal{F}$, 30.6 \mathcal{R} and 22 FPS of ReferFormer. We will then point-to-point analyze reasons of
340 improvements on $\mathcal{J}\&\mathcal{F}$ (for positive video), \mathcal{R} (for negative videos) and FPS (for inference speed).

- 341 • $\mathcal{J}\&\mathcal{F}$: (1) We introduce the early-grounding module which employs both pixel-wise and
342 channel-wise attention to enable multimodal interaction. Different from the CM-FPN used in
343 ReferFormer that solely fuses features from text to video in pixel-level, our early-grounding
344 module first enables pixel-level bi-directional fusion and then generates dynamic kernels
345 using the fused text feature g' to modulate the video feature f' . The dynamic convolution
346 (channel-wise attention) is commonly used to decode dense masks from visual features and
347 is suitable to suppress irrelevant features. By equipping text-guided dynamic convolution in
348 early-stage, the pixel decoder can be more focused on the target object (as shown in Figure 4).
349 (2) Our method leverages relational cycle consistency to constraint the intermediate feature
350 f_{early} to contain correct object-level information to recover some properties of original text
351 embedding. By applying this constraint, our method can better avoid interference and easier
352 locate the correct object. (3) Our instance query is composed of both original sentence
353 embedding and the reconstructed one. Different from ReferFormer that only utilizes original
354 sentence embedding as queries, the reconstructed embedding can encode visual information
355 to facilitate the instance query decode the objects from visual features.
- 356 • \mathcal{R} : The newly introduced metric \mathcal{R} aims to measure the robustness of the model against
357 unpaired inputs. Text-video pairs with (object-level) semantic consensus can be assumed
358 as in-distribution for RVOS problem where semantic consensus can be kind of easily
359 modeled. In contrast, unpaired text-video is much more difficult to tackle because there can
360 be unlimited out-of-distribution (OOD) scenarios for the text-video pairs. In our method,
361 instead of directly detect the OOD of input pairs, we convert the problem to find semantic
362 alignment between the input text embedding and reconstructed embedding and constraint
363 the property of reconstructed space by introducing the cycle consistency. In this way,
364 the comparison is conducted in the constraint original and reconstructed text spaces. For
365 ReferFormer, it models the alignment of text to video by querying the visual features by text
366 in the transformer decoder. In this way, the comparison is conducted in unconstrained text
367 and video spaces thus results in a inferior performance.
- 368 • FPS: The speed improvement of our method mainly comes from our efficient multimodal
369 fusion. Compared to the multi-scale CM-FPN, our early-grounding module is only conduct
370 at the high-level. In addition, our bi-direction multimodal fusion (Equ 4 & 5) only leverages
371 cross-attention to avoid computational expensive video-to-video operations.

372 A.2 Limitations

373 An important challenge for video segmentation is that target object disappearance due to occlusion,
374 which can results in false positives on a per-frame level. In our method, we predict the video-level
375 semantic alignment to handle the false positive in video-level resulted from unpaired text-video pairs.
376 However, since only video-level object expression is available in RVOS task, our method can not
377 address the frame-level false positives resulted from occlusion.

378 A.3 Additional Experiment on Negative Videos without Positive Text

Negative Video Source	\mathcal{R}	
	ReferFormer	Ours
Ref-Youtube-VOS	30.6	94.1
Ref-Youtube-VOS & Ref-DAVIS	33.1	92.2

Table A: Impact of different negative video sources.

379 As shown in Table A, we test the robustness of our model on two settings. We generate negative
380 videos from Ref-Youtube-VOS and a combination of Ref-Youtube-VOS and Ref-DAVIS dataset.

381 In both settings, all videos in the validation set are leveraged. The results indicates that source of
382 negative videos has minor impact on the robustness of our model.

383 **References**

- 384 [1] Adam Botach, Evgenii Zheltonozhskii, and Chaim Baskin. End-to-end referring video object
385 segmentation with multimodal transformers. *arXiv preprint arXiv:2111.14821*, 2021.
- 386 [2] Kan Chen, Jiyang Gao, and Ram Nevatia. Knowledge aided consistency for weakly supervised
387 phrase grounding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern
388 Recognition*, pages 4042–4050, 2018.
- 389 [3] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng,
390 and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference
391 on computer vision*, pages 104–120. Springer, 2020.
- 392 [4] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu.
393 Dynamic convolution: Attention over convolution kernels. In *Proceedings of the IEEE/CVF
394 Conference on Computer Vision and Pattern Recognition*, pages 11030–11039, 2020.
- 395 [5] Quan Cui, Boyan Zhou, Yu Guo, Weidong Yin, Hao Wu, and Osamu Yoshie. Zerovl: A strong
396 baseline for aligning vision-language representations with limited resources. *arXiv preprint
397 arXiv:2112.09331*, 2021.
- 398 [6] Zihan Ding, Tianrui Hui, Shaofei Huang, Si Liu, Xuan Luo, Junshi Huang, and Xiaoming Wei.
399 Progressive multimodal interaction network for referring video object segmentation. *The 3rd
400 Large-scale Video Object Segmentation Challenge*, page 7, 2021.
- 401 [7] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus
402 Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual
403 grounding. *arXiv preprint arXiv:1606.01847*, 2016.
- 404 [8] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and
405 Hongsheng Li. Dynamic fusion with intra-and inter-modality attention flow for visual question
406 answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
407 recognition*, pages 6639–6648, 2019.
- 408 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
409 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
410 pages 770–778, 2016.
- 411 [10] Drew A Hudson and Christopher D Manning. Compositional attention networks for machine
412 reasoning. *arXiv preprint arXiv:1803.03067*, 2018.
- 413 [11] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas
414 Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings
415 of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790, 2021.
- 416 [12] Guolin Ke, Di He, and Tie-Yan Liu. Rethinking positional encoding in language pre-training.
417 *arXiv preprint arXiv:2006.15595*, 2020.
- 418 [13] Anna Khoreva, Anna Rohrbach, and Bernt Schiele. Video object segmentation with language
419 referring expressions. In *Asian Conference on Computer Vision*, pages 123–141. Springer, 2018.
- 420 [14] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. *Advances in
421 neural information processing systems*, 31, 2018.
- 422 [15] Jin-Hwa Kim, Kyoung-Woon On, Woosang Lim, Jeonghee Kim, Jung-Woo Ha, and Byoung-
423 Tak Zhang. Hadamard product for low-rank bilinear pooling. *arXiv preprint arXiv:1610.04325*,
424 2016.
- 425 [16] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics
426 quarterly*, 2(1-2):83–97, 1955.
- 427 [17] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. Less
428 is more: Clipbert for video-and-language learning via sparse sampling. In *Proceedings of the
429 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7331–7341, 2021.
- 430 [18] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. Tvqa: Localized, compositional video
431 question answering. *arXiv preprint arXiv:1809.01696*, 2018.
- 432 [19] Muchen Li and Leonid Sigal. Referring transformer: A one-step approach to multi-task visual
433 grounding. *Advances in Neural Information Processing Systems*, 34, 2021.

- 434 [20] Xiaoya Li, Xiaofei Sun, Yuxian Meng, Junjun Liang, Fei Wu, and Jiwei Li. Dice loss for
435 data-imbalanced nlp tasks. *arXiv preprint arXiv:1911.02855*, 2019.
- 436 [21] Chen Liang, Yu Wu, Yawei Luo, and Yi Yang. Clawcranenet: Leveraging object-level relation
437 for text-based video segmentation. *arXiv preprint arXiv:2103.10702*, 2021.
- 438 [22] Chen Liang, Yu Wu, Tianfei Zhou, Wenguan Wang, Zongxin Yang, Yunchao Wei, and Yi Yang.
439 Rethinking cross-modal interaction from a top-down perspective for referring video object
440 segmentation. *arXiv preprint arXiv:2106.01061*, 2021.
- 441 [23] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie.
442 Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on
443 computer vision and pattern recognition*, pages 2117–2125, 2017.
- 444 [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense
445 object detection. In *Proceedings of the IEEE international conference on computer vision*,
446 pages 2980–2988, 2017.
- 447 [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
448 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European
449 conference on computer vision*, pages 740–755. Springer, 2014.
- 450 [26] Yang Liu, Samuel Albanie, Arsha Nagrani, and Andrew Zisserman. Use what you have: Video
451 retrieval using representations from collaborative experts. *arXiv preprint arXiv:1907.13487*,
452 2019.
- 453 [27] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
454 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
455 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 456 [28] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin
457 transformer. *arXiv preprint arXiv:2106.13230*, 2021.
- 458 [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint
459 arXiv:1711.05101*, 2017.
- 460 [30] Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon
461 Bharti, and Ming Zhou. Univl: A unified video and language pre-training model for multimodal
462 understanding and generation. *arXiv preprint arXiv:2002.06353*, 2020.
- 463 [31] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin
464 Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings
465 of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016.
- 466 [32] Antoine Miech, Ivan Laptev, and Josef Sivic. Learning a text-video embedding from incomplete
467 and heterogeneous data. *arXiv preprint arXiv:1804.02516*, 2018.
- 468 [33] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In
469 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
470 3967–3976, 2019.
- 471 [34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
472 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
473 style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*, 2019.
- 474 [35] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and
475 Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video
476 object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern
477 recognition*, pages 724–732, 2016.
- 478 [36] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung,
479 and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint
480 arXiv:1704.00675*, 2017.
- 481 [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
482 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
483 models from natural language supervision. In *International Conference on Machine Learning*,
484 pages 8748–8763. PMLR, 2021.

- 485 [38] Hamid Rezaatofghi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio
486 Savarese. Generalized intersection over union: A metric and a loss for bounding box regression.
487 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
488 pages 658–666, 2019.
- 489 [39] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object
490 segmentation network with a large-scale benchmark. In *Computer Vision–ECCV 2020: 16th*
491 *European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, pages
492 208–223. Springer, 2020.
- 493 [40] Xiaomeng Song, Yucheng Shi, Xin Chen, and Yahong Han. Explore multi-step reasoning
494 in video question answering. In *Proceedings of the 26th ACM international conference on*
495 *Multimedia*, pages 239–247, 2018.
- 496 [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
497 undefinedukasz Kaiser, and Illia Polosukhin. Attention is all you need. NIPS’17, page
498 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc.
- 499 [42] Liwei Wang, Yin Li, and Svetlana Lazebnik. Learning deep structure-preserving image-text
500 embeddings. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
501 pages 5005–5013, 2016.
- 502 [43] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and
503 Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the*
504 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8741–8750, 2021.
- 505 [44] Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring
506 video object segmentation. *arXiv preprint arXiv:2201.00487*, 2022.
- 507 [45] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas
508 Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint*
509 *arXiv:1809.03327*, 2018.
- 510 [46] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *Proceedings of the*
511 *IEEE/CVF International Conference on Computer Vision*, pages 5188–5197, 2019.
- 512 [47] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. Stacked attention networks
513 for image question answering. In *Proceedings of the IEEE conference on computer vision and*
514 *pattern recognition*, pages 21–29, 2016.
- 515 [48] Linwei Ye, Mrigank Roohan, Zhi Liu, and Yang Wang. Cross-modal self-attention network
516 for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer*
517 *Vision and Pattern Recognition*, pages 10502–10511, 2019.
- 518 [49] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling
519 context in referring expressions. In *European Conference on Computer Vision*, pages 69–85.
520 Springer, 2016.
- 521 [50] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks
522 for visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision*
523 *and pattern recognition*, pages 6281–6290, 2019.
- 524 [51] Zhou Yu, Jun Yu, Jianping Fan, and Dacheng Tao. Multi-modal factorized bilinear pooling with
525 co-attention learning for visual question answering. In *Proceedings of the IEEE international*
526 *conference on computer vision*, pages 1821–1830, 2017.
- 527 [52] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr:
528 Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*,
529 2020.

530 **Checklist**

- 531 1. For all authors...
- 532 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
533 contributions and scope? [Yes]
- 534 (b) Have you read the ethics review guidelines and ensured that your paper conforms to
535 them? [Yes]
- 536 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 537 (d) Did you describe the limitations of your work? [Yes]
- 538 2. If you are including theoretical results...
- 539 We are not including theoretical results.
- 540 3. If you ran experiments...
- 541 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
542 mental results (either in the supplemental material or as a URL)? [No] The code are
543 proprietary; most of our used dataset are public available; we plan to release code and
544 novel dataset upon acceptance.
- 545 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
546 were chosen)? [Yes]
- 547 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
548 ments multiple times)? [No] We conduct several experiments multiple times and find
549 the results are very close.
- 550 (d) Did you include the total amount of compute and the type of resources used (e.g., type
551 of GPUs, internal cluster, or cloud provider)? [Yes]
- 552 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 553 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 554 (b) Did you mention the license of the assets? [Yes]
- 555 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 556 (d) Did you discuss whether and how consent was obtained from people whose data you're
557 using/curating? [Yes] The data we used are public available datasets for academic
558 purposes.
- 559 (e) Did you discuss whether the data you are using/curating contains personally identifiable
560 information or offensive content? [No] The dataset that we used does not contain such
561 content.
- 562 5. If you used crowdsourcing or conducted research with human subjects...
- 563 We do not use crowdsourcing or conducted research with human subjects.