

# OnlineRefer: A Simple Online Baseline for Referring Video Object Segmentation

## Supplementary Material

Dongming Wu<sup>1‡</sup>, Tiancai Wang<sup>2</sup>, Yuang Zhang<sup>3</sup>, Xiangyu Zhang<sup>2,4</sup>, Jianbing Shen<sup>5†</sup>  
<sup>1</sup> Beijing Institute of Technology, <sup>2</sup> MEGVII Technology, <sup>3</sup> Shanghai Jiao Tong University,  
<sup>4</sup> Beijing Academy of Artificial Intelligence, <sup>5</sup> SKL-IOTSC, CIS, University of Macau

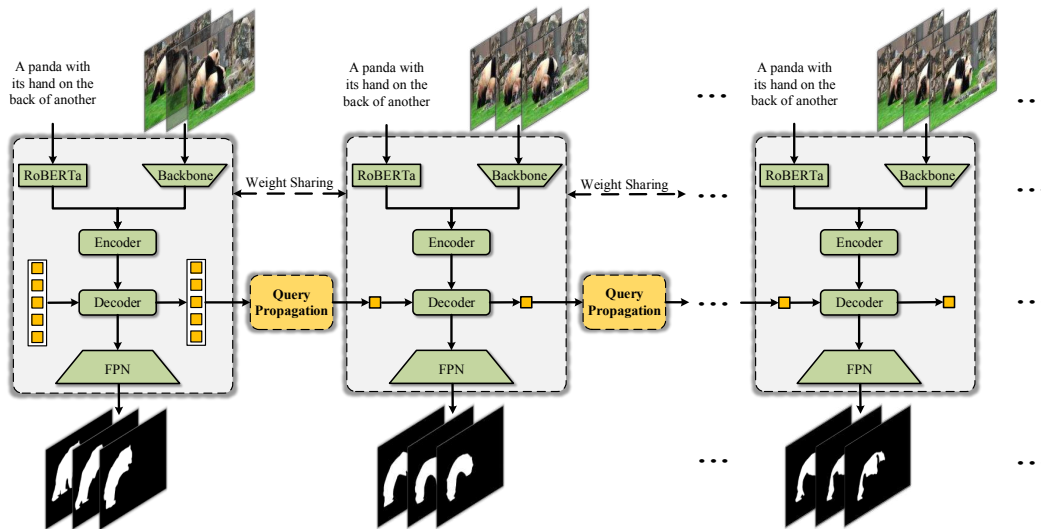


Figure 1: Pipeline of semi-online model of OnlineRefer, which associates the same referent object across different clips.

### 1. Semi-Online Model

We present more semi-online framework details of OnlineRefer in Fig. 1. Unlike the online model that follows a frame-by-frame pattern, the semi-online model propagates the target query across clips. Note that we use sharing query for multi-frame referring segmentation within each clip.

### 2. Additional Experiment Details

**Refer-Youtube-VOS** provides full-video expression by describing an entire video and first-frame expression based on the first frame, while we only use their full-video expression for training and validation.

**Refer-DAVIS<sub>17</sub>** also contains the full-video and first-frame expressions, which are developed by four annotators. Our final  $\mathcal{J}\&\mathcal{F}$  scores are averaged from the four results.

### 3. Additional Ablation Study

We provide a thorough analysis of sampling length settings in Table 1. It is obvious that increasing frames from 2 to 3 brings performance improvement on ResNet-50, while it fails on Swin-L. When using the progressive sampling

strategy (*i.e.*, [2, 3]), the best results can be obtained on both two backbones. This indicates that the appropriate increment on sampling lengths is beneficial for guaranteeing model stability and improving model performance.

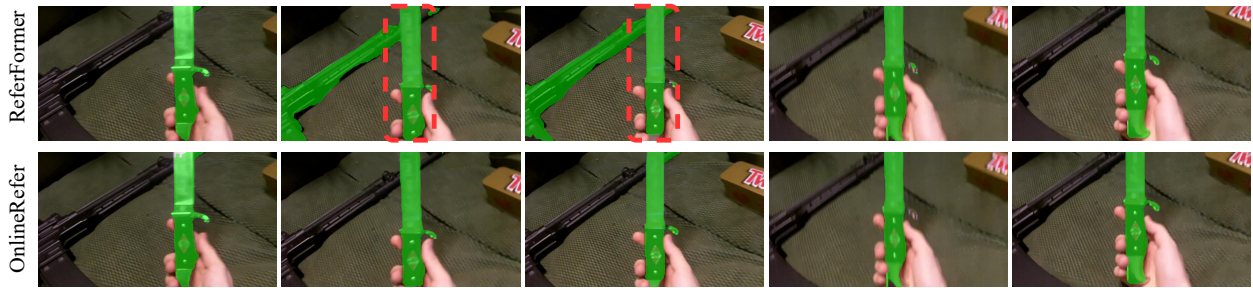
Sampling Lengths	ResNet-50			Swin-L		
	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$
[2]	56.5	55.1	57.9	62.7	60.8	64.5
[3]	56.7	55.2	58.2	/	/	/
[2, 3]	<b>57.3</b>	<b>55.8</b>	<b>58.8</b>	<b>63.5</b>	<b>61.6</b>	<b>65.5</b>
[2, 3, 4]	57.0	55.5	58.4	63.1	61.1	65.1

Table 1: The effect of sampling lengths on Refer-Youtube-VOS. ‘/’ means no results due to model divergence.

### 4. More Qualitative Results

Fig. 2 offers some qualitative comparison between the offline method ReferFormer [1] and our OnlineRefer. We can see that OnlineRefer performs better than ReferFormer under the situations of object occlusion and visually-similar background, approving the superiority of query propagation. Fig. 3 also shows OnlineRefer can deal well with other challenging situations, like size and appearance variation, small or missing objects, moving objects, *etc.*

Expression: a knife is in the hand of a person



Expression: the black and white zebra is on the left in the grass with its head down



Expression: a person in a blue shirt and black shorts



Expression: a red cloth being held by a person wearing black pants



Expression: a hedgehog next to a purple igloo

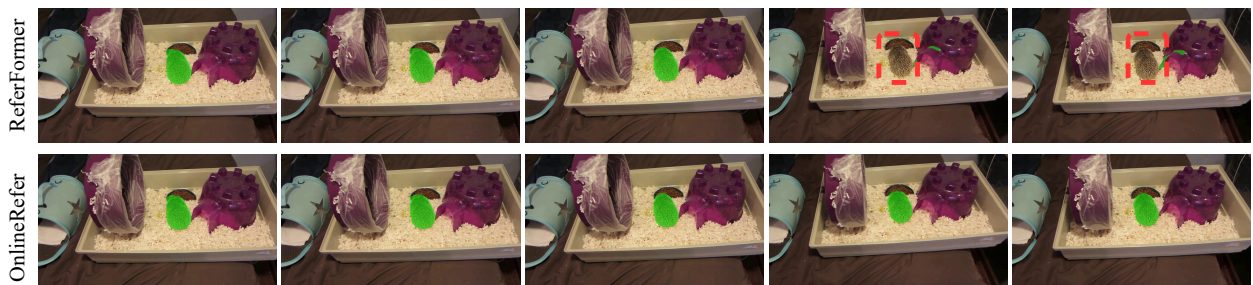


Figure 2: Comparison between ReferFormer and OnlineRefer on Refer-Youtube-VOS.

Expression: a man behind another man in a harness



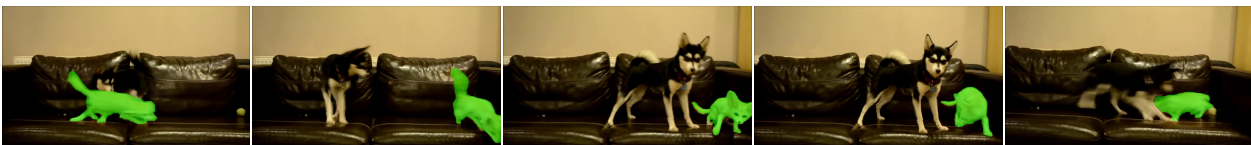
Expression: a large train racing down the tracks



Expression: a tennis racket on the left being held by a person



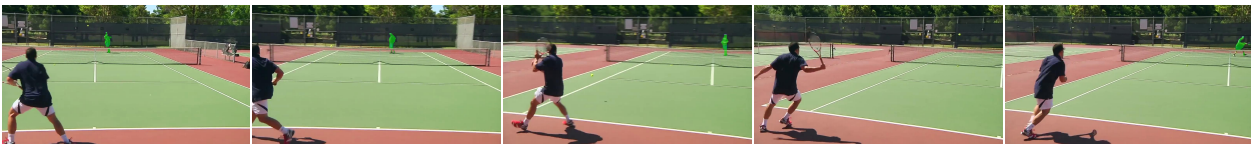
Expression: a small fox like dog on the right side of a sofa



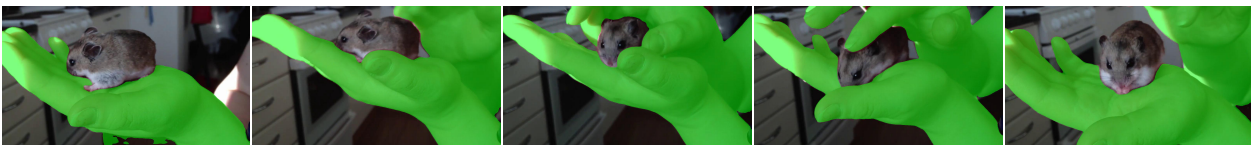
Expression: the second giraffe from the right



Expression: a person on the far side of a tennis court serving a tennis ball



Expression: the palm of a person carrying a mouse



Expression: a black cat in the middle of the view

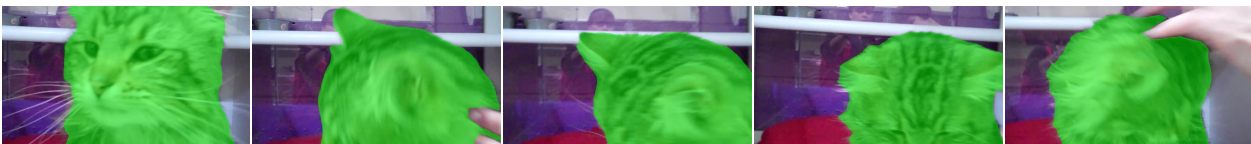


Figure 3: **More qualitative results** of our OnlineRefer on Refer-Youtube-VOS.

## References

- [1] Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring video object segmentation.

In *CVPR*, 2022.