# SeqFormer: Sequential Transformer for Video Instance Segmentation

Junfeng Wu[1*], Yi Jiang[2], Song Bai[2], Wenqing Zhang[1*], and Xiang Bai[1†]

[1] Huazhong University of Science and Technology [2] Bytedance

**Abstract.** In this work, we present SeqFormer for video instance segmentation. SeqFormer follows the principle of vision transformer that models instance relationships among video frames. Nevertheless, we observe that a stand-alone instance query suffices for capturing a time sequence of instances in a video, but attention mechanisms shall be done with each frame independently. To achieve this, SeqFormer locates an instance in each frame and aggregates temporal information to learn a powerful representation of a video-level instance, which is used to predict the mask sequences on each frame dynamically. Instance tracking is achieved naturally without tracking branches or post-processing. On YouTube-VIS, SeqFormer achieves 47.4 AP with a ResNet-50 backbone and 49.0 AP with a ResNet-101 backbone without bells and whistles. Such achievement significantly exceeds the previous state-of-the-art performance by 4.6 and 4.4, respectively. In addition, integrated with the recently-proposed Swin transformer, SeqFormer achieves a much higher AP of 59.3. We hope SeqFormer could be a strong baseline that fosters future research in video instance segmentation, and in the meantime, advances this field with a more robust, accurate, neat model. The code is available at https://github.com/wjf5203/SeqFormer.

**Keywords:** Video instance segmentation, Video transformer

## 1 Introduction

Video Instance Segmentation (VIS) [45,30] is an emerging vision task that aims to simultaneously perform detection, classification, segmentation, and tracking of object instances in videos. Compared to image instance segmentation [11], video instance segmentation is much more challenging since it requires accurate tracking of objects across an entire video.

Previous VIS algorithms can be roughly divided into two categories. The first mainstream follows the tracking-by-detection paradigm, extending image instance segmentation models with a tracking branch [45,4,46,21,28,10]. These methods first predict candidate detection and segmentation frame-by-frame, and then associate them by classification [45,46] or re-identification [4,21] to track

---

∗ Work done during an internship at ByteDance.
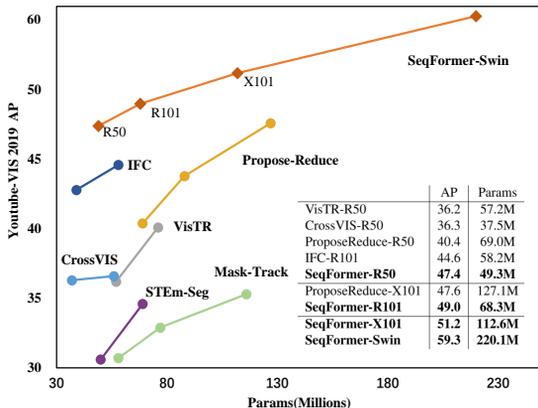† Corresponding author: xbai@hust.edu.cn

| | AP | Params |
|---|---|---|
| VisTR-R50 | 36.2 | 57.2M |
| CrossVIS-R50 | 36.3 | 37.5M |
| ProposeReduce-R50 | 40.4 | 69.0M |
| IFC-R101 | 44.6 | 58.2M |
| **SeqFormer-R50** | **47.4** | **49.3M** |
| ProposeReduce-X101 | 47.6 | 127.1M |
| **SeqFormer-R101** | **49.0** | **68.3M** |
| **SeqFormer-X101** | **51.2** | **112.6M** |
| **SeqFormer-Swin** | **59.3** | **220.1M** |

**Fig. 1. Performance *vs.* Model Size.** All results are reported with single model and single-scale inference. SeqFormer significantly outperforms the previous method with similar parameters.

the instance through a video. Nevertheless, the tracking process is sensitive to occlusions and motion blur that are common in videos. Another mainstream is to predict clip-level instance masks by taking a video clip [2,3] or the entire video [14,39] as input. It divides a video into multiple overlapping clips and generates mask sequences with clip-by-clip matching on overlapping frames. More recently, VisTR [39] first adapts transformer [37] to VIS and uses instance queries to obtain instance sequence from video clips. After that, IFC [14] improves the performance and efficiency of VisTR by building communications between frames in a transformer encoder.

In this paper, we present Sequential Transformer (SeqFormer), which follows the principle of vision transformer [5,39] and models instance relationships among video frames. As in [14], we observe that a stand-alone instance query suffices although an object may be of different positions, sizes, shapes, and various appearances. Nevertheless, it is witnessed that the attention process shall be done with each frame independently, so that the model will attend to locations following with the movement of instance through the video. This observation aligns with the conclusion drawn in action recognition [47,29], where the 1D time domain and 2D space domain have different characteristics and should be handled in a different fashion.

Considering the movement of an instance in a video, a model is supposed to attend to different spatial locations following the motion of the instance. We decompose the shared instance query into frame-level box queries for the attention mechanism to guarantee that the attention focuses on the same instance on each frame. The box queries are kept on each frame and used to predict the bounding box sequences. Then the features within the bounding boxes are aggregated to refine the box queries on the current frame. By repeating this refinement through

decoder layers, SeqFormer locates the instance in each frame in a coarse-to-fine manner, in a similar way to Deformable DETR [49].

However, to mitigate redundant information from non-instance frames, those box queries are aggregated in a weighted manner, where the weights are end to end learned upon the box embeddings. The generated representation, which retains richer object cues, is used to predict the category and generate dynamic convolution weights of mask head. Since the box sequences are predicted and refined in the decoder, SeqFormer naturally and succinctly establishes the association of instances across frames.

In summary, SeqFormer enjoys the following advantages:

- SeqFormer is a neat and efficient end-to-end framework. Given an arbitrary long video as input, SeqFormer predicts the classification results, box sequences, and mask sequences in one step without the need for additional tracking branches or hand-craft post-processing.
- As shown in Fig. 1, SeqFormer sets the new state-of-the-art performance on YouTube-VIS 2019 benchmark [45]. SeqFormer achieves 47.4 AP with a ResNet-50 backbone and 49.0 AP with a ResNet-101 backbone without bells and whistles. Such achievement significantly exceeds the previous state-of-the-art performance by 4.6 and 4.4, respectively. With a ResNext-101 backbone, SeqFormer achieves 51.2 AP, which is the first time that an algorithm achieves an AP above 50. In addition, integrated with the recently-proposed Swin transformer, SeqFormer achieves a much higher AP of 59.3.
- With the query decomposition mechanism, SeqFormer attends to locations following with the movement of instance through the video and learns a powerful representation for instance sequences.
- The code and the pre-trained models are publicly available. We hope the SeqFormer, with the idea of making attention follow with the movement of object, could be a strong baseline that fosters future research in video instance segmentation, and in the meantime, advances this field with a more robust, accurate, neat model.

## 2    Related Work

**Image Instance Segmentation** Instance Segmentation is the most fundamental and challenging task in computer vision, which aims to detect every instance and segment every pixel respectively in static images. Instance segmentation was dominated by Mask R-CNN architecture [11,22,13] for a long time, Mask R-CNN [11] directly introduces fully convolutional mask head to Faster R-CNN [31] in a multi-task learning manner. Recently, one stage models [6,38,41,35] emerged as excellent frameworks for instance segmentation. Solo [38] and CondInst [35] propose one stage instance segmentation pipeline and achieve comparable performance. CondInst [35] proposes to dynamically generate the mask head parameters for each instance, which is used to predict the mask of the corresponding instance. QueryInst [8] proposes a query based instance segmentation framework based on Sparse R-CNN [34], which also take advantage of the Dynamic mask

head. Dynamic mask head can be efficiently adopted into video segmentation tasks because instances with the same identity on different frames can share the same mask head parameters.

**Video Instance Segmentation.** Video instance segmentation is extended from the traditional image instance segmentation, and aims to simultaneously segment and track all object instances in the video sequence. The baseline method Mask-Track R-CNN [45] is built upon Mask R-CNN [11] and introduces a tracking head to associate each instance in the video. SipMask [4] proposes a spatial preservation module to generate spatial coefficients for mask predictions based on the one-stage FCOS [36]. STMask [17] proposes a spatial feature calibration to extract features for frame-level instance segmentation on each frame, and further introduces a temporal fusion module to aggregate temporal information from adjacent frames. STEm-Seg [2] models a video clip as a single 3D spatial-temporal volume and enables inference procedure based on clustering. CrossVIS [46] proposes a learning scheme that uses the instance feature in the current frame to pixel-wisely localize the same instance in other frames. MaskProp [3] and Propose-Reduce [18] take advantage of mask propagation, which can achieve high performance, but it is very computationally intensive.

**Transformers.** Transformer [37] was first proposed for the sequence-to-sequence machine translation task and became the basic component in most Natural Language Processing tasks. Recently, Transformer has been successfully applied in many visual tasks such as Object detection [5,49,34], segmentation [42,48,39,40], tracking [26,33,44], video recognition [29,16,1,23]. VIT[7] firstly applies transformer in image recognition and model an image as sequence of patches, which achieves comparable performance with traditional CNN architecture. DETR [5] proposes a new detection paradigm upon transformers, which simplifies the traditional detection framework and abandons the hand-crafted post-processing module. Deformable DETR [49] achieves better performance by using local attention and multi-scale feature maps. VisTR [39] is the first method that adapts Transformer to the VIS task. However, VisTR has a fixed number of input queries hardcoded by video length and maximum number of instances. Each query corresponds to an object on every single frame. In our method, instances with the same identity share a same query, which aggregates information across the video and learn a global feature representation for efficient segmentation. IFC [14] improves the performance and efficiency of VisTR by building communications between frames in the transformer encoder instead of flatting the space-time features into one dimension, but it still flatten the space-time features for the transformer decoder. Our model is designed to carry out the instance feature capturing independently on different frames, which makes the model attend to locations following with the movement of instance through the video.
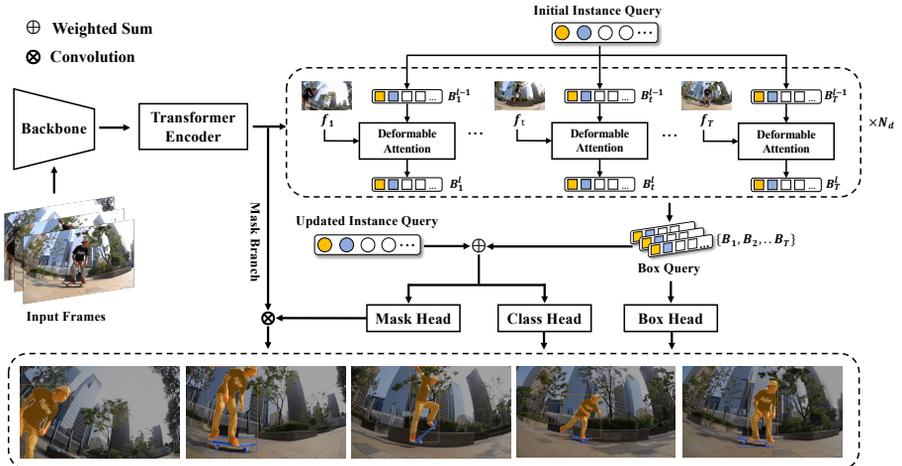
**Fig. 2.** The overall architecture of SeqFormer. Given the feature maps of input frames, the initial instance query is decomposed into frame-level box queries at the first decoder layer. The box queries are kept on each frame and serve as anchors without interacting with each other. The features extracted by box queries from each frame are aggregated to the instance query after each decoder layer, which is used for predicting dynamic mask head parameters. Then the mask head convolves the encoded feature maps to generate the mask sequences.

## 3    Method

### 3.1    Architecture

The network architecture is visualized in Fig. 2. SeqFormer has a CNN backbone and a transformer encoder for extracting feature maps from each frame independently. Next, a transformer decoder is adapted to locate the instance sequences and generate a video-level instance representation. Finally, three output heads are used for instance classification, instance sequences segmentation, and bounding box prediction, respectively.

**Backbone.** Given an input video $x_v \in \mathbb{R}^{T \times 3 \times H \times W}$ with 3 color channels and $T$ frames of resolution $H \times W$, the CNN backbone (*e.g.*, ResNet [12]) extracts feature maps for each frame independently.

**Transformer Encoder.** First, a $1 \times 1$ convolution is used to reduce the channel dimension of the all the feature maps to $C = 256$, creating new feature maps $\{\mathbf{f}'_t\}_{t=1}^T, \mathbf{f}'_t \in \mathbb{R}^{C \times H' \times W'}, t \in [1, T]$. After adding fixed positional encodings [5], the transformer encoder performs deformable attention [49] on the feature maps, resulting in the output feature maps $\{\mathbf{f}_t\}_{t=1}^T$, with the same resolutions as the input. To perform attention mechanisms on each frame independently, we keep the spatial and temporal dimensions of feature maps rather than flattening them into one dimension.

**Query Decompose Transformer Decoder.** Given a video, humans can effortlessly identify every instance and associate them through the video, despite the various appearance and changing positions on different frames. If an instance is hard to recognize due to occlusion or motion blur in some frames, humans can still re-identify it through the context information from other frames. In other words, for the same instance on different frames, humans treat them as a whole instead of individuals. This is the crucial difference between video and image instance segmentation. Motivated by this, we propose Query Decompose Transformer Decoder, which aims to learn a more and robust video-level instance representation across frames.

We introduce a fixed number of learnable embeddings to query the features of the same instance from each frame, termed Instance Queries. Different from the instance queries corresponding to frame-level instances in VisTR [39], which has a fixed number of input queries hardcoded by video length and maximum number of instances, our instance queries correspond to video-level instances. Since the changing appearance and position of the instance, the model should attend to different exact spatial locations of each frame. To achieve this goal, we propose to decompose the instance query into $T$ frame-specific box queries, each of which serves as an anchor for retrieving and locating features on the corresponding frame.

At the first decoder layer, an instance query $\mathbf{I}_q \in \mathbb{R}^C$ is used to query the instance features on features maps of each frame independently:

$$\mathbf{B}_t^1 = \text{DeformAttn}(\mathbf{I}_q, \mathbf{f}_t), \tag{1}$$

where $\mathbf{B}_t^1 \in \mathbb{R}^C$ is the box query on frame $t$ from the 1-st decoder layer, and DeformAttn indicates deformable attention module in [49]. Given a query element and the frame feature map $\mathbf{f}_t$, deformable attention only attends to a small set of key sampling points. At the $l$-th ($l > 1$) layer, the box query $\mathbf{B}_t^{l-1}$ from last layer is used as input:

$$\mathbf{B}_t^l = \text{DeformAttn}(\mathbf{B}_t^{l-1}, \mathbf{f}_t), \tag{2}$$

and the instance query aggregates the temporal features by a weighted sum of all the box queries at the end of every decoder layers, where the weights are end to end learned upon the box embedding:

$$\mathbf{I}_q^l = \frac{\sum_{t=1}^T \mathbf{B}_t^l \times \text{FC}(\mathbf{B}_t^l)}{\sum_{t=1}^T \text{FC}(\mathbf{B}_t^l)} + \mathbf{I}_q^{l-1}. \tag{3}$$

After $N_d$ decoder layers, we get an instance query and $T$ box queries for each instance. The instance query is a shared video-level instance representation, and the box query contains the position information for predicting the bound box on each frame. We define the instance query $I_q^{N_d}$ and box queries $\{B_t^{N_d}\}_{t=1}^T$ from the last layer of decoder as output instance embedding and box embeddings $\{\mathbf{BE}_t\}_{t=1}^T, \mathbf{BE}_t \in \mathbb{R}^{N \times d}$.

**Output Heads.** As shown in Fig. 2, we add mask head, box head, class head on the top of the decoder outputs. A linear projection acts as the class head

to produce the classification results. Given the instance embedding from the transformer decoder with index $\sigma(i)$, class head output a class probability of class $\mathbf{c}_i$ (which may be $\varnothing$) as $\hat{p}_{\sigma(i)}(\mathbf{c}_i)$ .

The box head is a 3-layer feed forward network (FFN) with ReLU activation function and a linear projection layer. For $\mathbf{BE}_t$ of each frame, the FFN predicts the normalized center coordinates, height and width of the box w.r.t. the frame. Thus, for the instance with index $\sigma(i)$, we denote the predicted box sequence as $\hat{\mathbf{b}}_{\sigma(i)} = \{\hat{\mathbf{b}}_{(\sigma(i),1)}, \hat{\mathbf{b}}_{(\sigma(i),2)}, ..., \hat{\mathbf{b}}_{(\sigma(i),T)}\}$.

As for mask head, we leverage dynamic convolution [35] as our mask head. The output instance embedding of decoder contains the information of instance on all frames, thus it can be regarded as a more robust instance representation. We can use instance embedding to efficiently generate the entire mask sequences. First, a 3-layer FFN encodes the instance embedding into parameters $\omega_i$ of mask head with index $\sigma(i)$, which has three $1 \times 1$ convolution layers. The instances with the same identity on different frames share the same mask head parameters, which makes the segmentation very efficient. Each convolution layer has 8 channels and uses ReLU as the activation function except for the last one, following [35]. As shown in Fig. 2, there is a mask branch that provides the feature maps for mask head to predict instance masks. We employ an FPN-like architecture as the mask branch to make use of multi-scale feature maps from transformer encoder and generate feature maps sequences $\{\hat{\mathbf{F}}^1_{\text{mask}}, \hat{\mathbf{F}}^2_{\text{mask}}, ..., \hat{\mathbf{F}}^T_{\text{mask}}\}$ that are $\frac{1}{8}$ of the input resolution and have 8 channels for each frame independently. Then the feature map $\hat{\mathbf{F}}^t_{\text{mask}}$ is concatenated with a map of the relative coordinates from center of $\hat{\mathbf{b}}_{(\sigma(i),t)}$ in corresponding frames to provide a location cue for predicting the instance mask. Thus we get the $\{\mathbf{F}^t_{\text{mask}}\}^T_{t=1}, \mathbf{F}^t_{\text{mask}} \in \mathbb{R}^{10 \times \frac{H}{8} \times \frac{W}{8}}$. The mask head performs convolution on these high-resolution sequence feature maps $\mathbf{F}^t_{\text{mask}}$ to predict the mask sequences:

$$\{\mathbf{m}^t_i\}^T_{t=1} = \{\text{MaskHead}(\mathbf{F}^t_{\text{mask}}, \omega_i)\}^T_{t=1}, \qquad (4)$$

where MaskHead performs three-layer $1 \times 1$ convolution on given feature maps with the kernels reshaped from $\omega$. By sharing the same mask head parameters for instances with the same identity on different frames, our method can efficiently perform instance segmentation on each frame. Similar to DETR [5], we add output heads and Hungarian loss after each decoder layer as an auxiliary loss to supervise the training stage.

## 3.2   Instance Sequences Matching and Loss

Our method predicts a fixed-size set of $N$ predictions in a single pass through the decoder, and $N$ is set to be significantly larger than the number of instances in a video. To train our network, we first need to find a bipartite graph matching between the prediction and the ground truth. Let $\mathbf{y}$ denotes the ground truth set of video-level instance, and $\hat{\mathbf{y}}_i = \{\hat{\mathbf{y}}_i\}^N_{i=1}$ denotes the predicted instance set. Each element $i$ of the ground truth set can be seen as $\mathbf{y}_i = \{\mathbf{c}_i, (\mathbf{b}_{i,1}, \mathbf{b}_{i,2}, ..., \mathbf{b}_{i,T})\}$,

where $\mathbf{c}_i$ is the target class label including $\varnothing$, and $\mathbf{b}_{i,t} \in [0,1]^4$ is a vector that defines ground truth bounding box center coordinates and its relative height and width in the frame $t$. For the predictions of instance with index $\sigma(i)$, we take the output of class head $\hat{p}_{\sigma(i)}(\mathbf{c}_i)$ and predicted bounding box $\hat{\mathbf{b}}_{\sigma(i)}$. Then we define the pair-wise matching cost between ground truth $y_i$ and a prediction with index $\sigma(i)$.

$$\mathcal{L}_{\text{match}}(\mathbf{y}_i, \hat{\mathbf{y}}_{\sigma(i)}) = -\hat{p}_{\sigma(i)}(\mathbf{c}_i) + \mathcal{L}_{\text{box}}(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)}), \qquad (5)$$

where $\mathbf{c}_i \neq \varnothing$. Note that Eq. 5 does not consider the similarity between mask prediction and mask ground truth, as such mask-level comparison is computationally expensive. To find the best assignment of a ground truth to a prediction, we search for a permutation of N elements $\sigma \in S_n$ with the lowest cost:

$$\hat{\sigma} = \underset{\sigma \in S_n}{\arg\min} \sum_i^N \mathcal{L}_{\text{match}}(\mathbf{y}_i, \hat{\mathbf{y}}_{\sigma(i)}). \qquad (6)$$

Following prior work[5,39], the optimal assignment is computed with the Hungarian algorithm [15]. Given the optimal assignment $\hat{\sigma}$, we use Hungarian loss for all matched pairs to train our network:

$$\mathcal{L}_{\text{Hung}}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{i=1}^N \Big[ -\log \hat{p}_{\hat{\sigma}(i)}(\mathbf{c}_i) + \mathbb{1}_{\{\mathbf{c}_i \neq \varnothing\}} \mathcal{L}_{\text{box}}(\mathbf{b}_i, \hat{\mathbf{b}}_{\hat{\sigma}}(i))$$
$$+ \mathbb{1}_{\{\mathbf{c}_i \neq \varnothing\}} \mathcal{L}_{\text{mask}}(\mathbf{m}_i, \hat{\mathbf{m}}_{\hat{\sigma}}(i)) \Big]. \qquad (7)$$

For $\mathcal{L}_{\text{box}}$, we use a linear combination of the $\mathcal{L}_1$ loss and the generalized IoU loss [32]. The mask sequences $\{m_i^t\}_{t=1}^T$ from mask head with $\frac{1}{8}$ of the video resolution which may loss some details, thus we upsample the predicted mask to $\frac{1}{4}$ of the video resolution, and downsample the ground truth mask to the same resolution for mask loss, following [35]. The mask loss $\mathcal{L}_{\text{mask}}$ is defined as a combination of the Dice [27] and Focal loss [19]. We calculate box loss and mask loss on each frame and take the average for Hungarian loss.

## 4    Experiment

### 4.1    Datasets and Metrics

We evaluate our method on YouTube-VIS 2019 [45] and YouTube-VIS 2021 [43] datasets. YouTube-VIS 2019 is the first and largest dataset for video instance segmentation, which contains 2238 training, 302 validation, and 343 test high-resolution YouTube video clips. It has a 40-category label set and 131k high-quality instance masks. In each video, objects with bounding boxes and masks are labeled every five frames. YouTube-VIS 2021 is an improved and extended version of YouTube-VIS 2019 dataset, it contains 3,859 high-resolution videos and 232k instance annotations. The newly added videos in the dataset include more instances and frames.

Video instance segmentation is evaluated by the metrics of average precision (AP) and average recall (AR). Different from image instance segmentation, each instance in a video contains a sequence of masks. To evaluate the spatio-temporal consistency of the predicted mask sequences, the IoU computation is carried out in the spatial-temporal domain. This requires a model not only to obtain accurate segmentation and classification results at frame-level but also to track instance masks between frames accurately.

## 4.2    Implementation Details

**Model settings.** ResNet-50 [12] is used as our backbone network unless otherwise specified. Similar to [49], we use the features from the last three stages as {C3, C4, C5} in ResNet, which correspond to the feature maps with strides {8, 16, 32}. And adding the lowest resolution feature map C6 obtained via a 3 × 3 stride 2 convolution on the C5. We set sampled key numbers K=4 and eight attention heads for deformable attention modules. We use six encoder and six decoder layers of hidden dimension 256 for the transformer, and the number of instance queries is set to 300.

**Training.** We used AdamW [25] optimizer with base learning rate of $2 \times 10^{-4}$, $\beta 1 = 0.9$ , $\beta 2 = 0.999$, and weight decay of $10^{-4}$. Learning rates of the backbone and linear projections used for deformable attention modules are multiplied by a factor of 0.1. We first pre-train the model on COCO [20] by setting the number of input frames $T = 1$. Given the pretrained weights, we train our models on the YouTube-VIS dataset with input frames $T = 5$ sampled from the same video.

The training data of the YouTube-VIS dataset is not sufficient, which makes a model prone to overfitting. To address this problem, we adopt 80K training images in the COCO for compensation, following [2,18]. We only use the images with 20 overlapping categories in COCO and augment them with ±10° rotation to generate a five-frame pseudo video. We train our model on the mixed dataset including COCO and the video dataset for 12 epochs, and the learning rate is decayed at the 6-th and 10-th epoch by a factor of 0.1. The input frame sizes are downsampled so that the longest side is at most 768 pixels. The model is implemented with PyTorch-1.7 and is trained on 8 V100 GPUs of 32G RAM, with 2 video clips per GPU.

**Inference.**   SeqFormer is able to model a video of arbitrary length without grouping frames into subsequences. We take the whole video as input during inference, which is downscaled to 360p, following MaskTrack R-CNN [45]. SeqFormer learns a video-level instance representation used for dynamic segmentation on each frame and classification, and the box sequences are generated by the decoder. Thus, no post-processing is needed for associating instances.

## 4.3    Main Results

The comparison of SeqFormer with previous state-of-the-art methods on YouTube-VIS 2019 are listed in Table 1. MaskProp [3] and ProposeReduce [18] are the state-of-the-art methods, which take a strong backbone to extract spatial

**Table 1.** Quantitative results of video instance segmentation on YouTube-VIS 2019 validation set. The result with superscript "†" is obtained without coco joint training. The best results with the same backbone are in **bold**.

| Backbone | Method | Params | FPS | AP | $AP_{50}$ | $AP_{75}$ | $AR_1$ | $AR_{10}$ |
|---|---|---|---|---|---|---|---|---|
| | MaskTrack R-CNN [45] | 58.1M | 20.0 | 30.3 | 51.1 | 32.6 | 31.0 | 35.5 |
| | STEm-Seg [2] | 50.5M | 7.0 | 30.6 | 50.7 | 33.5 | 37.6 | 37.1 |
| | SipMask [4] | 33.2M | 30.0 | 33.7 | 54.1 | 35.8 | 35.4 | 40.1 |
| | CompFeat [9] | - | - | 35.3 | 56.0 | 38.6 | 33.1 | 40.3 |
| | SG-Net [21] | - | - | 34.8 | 56.1 | 36.8 | 35.8 | 40.8 |
| ResNet-50 | VisTR [39] | 57.2M | 69.9 | 36.2 | 59.8 | 36.9 | 37.2 | 42.4 |
| | MaskProp [3] | - | - | 40.0 | - | 42.9 | - | - |
| | CrossVIS [46] | 37.5M | 39.8 | 36.3 | 56.8 | 38.9 | 35.6 | 40.7 |
| | Propose-Reduce [18] | 69.0M | - | 40.4 | 63.0 | 43.8 | 41.1 | 49.7 |
| | IFC [14] | 39.3M | 107.1 | 42.8 | 65.8 | 46.8 | 43.8 | 51.2 |
| | **SeqFormer**$^{†}$ | 49.3M | 72.3 | 45.1 | 66.9 | 50.5 | **45.6** | 54.6 |
| | **SeqFormer** | 49.3M | 72.3 | **47.4** | **69.8** | **51.8** | 45.5 | **54.8** |
| | MaskTrack R-CNN [45] | 77.2M | - | 31.8 | 53.0 | 33.6 | 33.2 | 37.6 |
| | STEm-Seg [2] | 69.6M | - | 34.6 | 55.8 | 37.9 | 34.4 | 41.6 |
| | SG-Net [21] | - | - | 36.3 | 57.1 | 39.6 | 35.9 | 43.0 |
| ResNet-101 | VisTR [39] | 76.3M | 57.7 | 40.1 | 64.0 | 45.0 | 38.3 | 44.9 |
| | MaskProp [3] | - | - | 42.5 | - | 45.6 | - | - |
| | CrossVIS [46] | 56.6 | 35.6 | 36.6 | 57.3 | 39.7 | 36.0 | 42.0 |
| | Propose-Reduce [18] | 88.1M | - | 43.8 | 65.5 | 47.4 | 43.0 | 53.2 |
| | IFC [14] | 58.3M | 89.4 | 44.6 | 69.2 | 49.5 | 44.0 | 52.1 |
| | **SeqFormer** | 68.4M | 64.6 | **49.0** | **71.1** | **55.7** | **46.8** | **56.9** |
| | MaskProp [3] | - | - | 44.3 | - | 48.3 | - | - |
| ResNeXt-101 | Propose-Reduce [18] | 127.1M | - | 47.6 | 71.6 | 51.8 | 46.3 | 56.0 |
| | **SeqFormer** | 112.7M | 30.7 | **51.2** | **75.3** | **58.0** | **46.5** | **57.3** |
| Swin-L | **SeqFormer** | 220.0M | 27.7 | **59.3** | **82.1** | **66.4** | **51.7** | **64.4** |

features and use mask propagation to improve the segmentation and tracking, but suffer from low inference speed. We list the methods with different backbones for fair comparison. It can be observed that SeqFormer significantly surpasses all the previous best reported results by at least 4 AP with the same backbone. Training our model with coco pseudo videos improves the AP from 45.1 to 47.4. SeqFormer with a ResNet-50 backbone can even achieve competitive performance against state-of-the-art methods with a ResNeXt-101 backbone. By adopting Swin transformer [24] as our backbone without further modifications, SeqFormer can first achieve 59.3 AP on this benchmark, outperforming the best previous results by a large margin of 11.7 AP. To understand the runtime efficiency, we measure FPS of SeqFormer excluding the data loading process of multiple images on NVIDIA Tesla V100. With an input size of 360p and a ResNet-50 backbone on YouTube-VIS, the inference FPS is 72.3. While surpassing the stateof-the-art AP by a large margin, SeqFormer is the second fast one following IFC. An example of qualitative comparison with previous meth-

**Table 2.** Quantitative results of video instance segmentation on YouTube-VIS 2021 validation set. The best results with the same backbone are in **bold**.

| Backbone | Method | AP | $AP_{50}$ | $AP_{75}$ | $AR_1$ | $AR_{10}$ |
|---|---|---|---|---|---|---|
| ResNet-50 | MaskTrack R-CNN [45] | 28.6 | 48.9 | 29.6 | 26.5 | 33.8 |
| | SipMask [4] | 31.7 | 52.5 | 34.0 | 30.8 | 37.8 |
| | CrossVIS [46] | 34.2 | 54.4 | 37.9 | 30.4 | 38.2 |
| | IFC [14] | 36.6 | 57.9 | 39.3 | - | - |
| | **SeqFormer** | **40.5** | **62.4** | **43.7** | **36.1** | **48.1** |
| Swin-L | **SeqFormer** | **51.8** | **74.6** | **58.2** | **42.8** | **58.1** |



**Fig. 3.** Qualitative comparisons with other methods.

ods is given in Fig. 3, the mask predictions of SeqFormer are more stable over time. Please refer to the Sup. Mat. for more qualitative results. We also evaluate our approach on the recently introduced YouTube-VIS 2021 dataset, which is a more challenging dataset with more videos and a higher number of instances and frames. As shown in Table 2, SeqFormer achieves 40.5 AP with a ResNet-50 backbone, surpassing previous methods by 3.9 AP. We believe that our effective method will serve as a strong baseline on these benchmarks and facilitate future research in video instance segmentation.

## 4.4   Ablation Study

This section conducts extensive ablation experiments to study the effects of different settings in our proposed method. All the ablation experiments are conducted with the ResNet-50 backbone and training on YouTube-VIS 2019 dataset rather than the mixed dataset.

**Table 3.** Instance query decomposition. Decomposing instance query into frame-level box queries is critical for SeqFormer.

| Decompose | AP | $AP_{50}$ | $AP_{75}$ | $AR_1$ | $AR_{10}$ |
|-----------|------|------|------|------|------|
| w/o       | 34.1 | 53.7 | 34.9 | 34.8 | 40.9 |
| w         | 45.1 | 66.9 | 50.5 | 45.6 | 54.6 |

**Instance query decomposition.** Instance query decomposition plays an important role in our method. Since an instance may have different positions on each frame, the iterative refinement of the spatial sampling region should be performed independently on each frame. To keep the temporal consistency of instances, we use the temporal-shared instance query for deformable attention and get box queries for each frame. The box queries will be kept through all the decoder layers and serve as frame anchors for the same instance. Experiments of models without box queries and using the shared instance query for each decoder layer are presented in Table 3. The model without query decomposition manages to achieve only 34.1 AP. It is because the query controls the sampling region of deformable attention. Using the same instance query for each frame will result in the same spatial sampling region on each frame, as shown in Fig. 4 (a), which is inaccurate and insufficient for video-level instance representation. We further visualize the sampling points of the second and the last decoder layers in Fig. 4 (b) and (c). The box queries decoupled from instance query serve as anchors for locating features and iteratively refining the sampling region on the current frame. It can be seen that SeqFormer attends to locations following with the movement of instance through the video in a coarse-to-fine manner. Please refer to the Sup. Mat. for more visualization of sampling points.

**Spatial temporal dimensions.** Previous transformer-based methods [39,14] flatten the spatial and temporal dimensions of video features into one dimension for the transformer decoder. We argue that the temporal dimension should not be flattened with spatial dimensions, since it was recognized that the 2D space domain and 1D time domain have different characteristics and should be intuitively handled in a different way [47]. Thus, we retain the complete 3D spatio-temporal dimensions and perform explicit region sampling and information aggregation on all frames. In this experiment, we study the effect of this architecture by replacing deformable transformer with vanilla transformer and flattening the spatial and temporal dimensions, termed as 'flatten' in Table 4. For fair comparison, we use single-scale deformable attention as the baseline, termed as 'single-scale', which use the same scale feature map with 'flatten', the default setting termed as 'multi-scale'. By keeping spatial-temporal dimensions of video features, the AP increased from 35.1 to 42.5. The use of multi-scale feature maps can only improve 2.6 AP, which proves that the success of our method mainly comes from the preservation of the temporal dimension and the explicit spatial sampling.
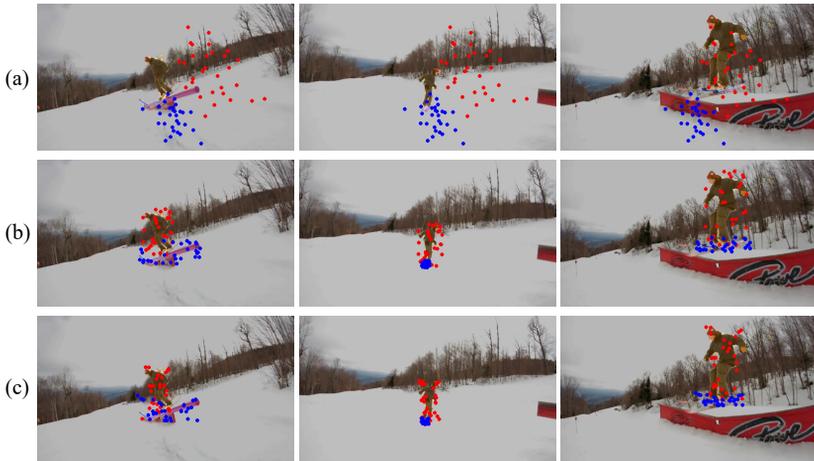
**Fig. 4.** The sampling points from the first decoder layer is shown in (a), which is coarse and inaccurate. The refined accurate sampling points from the second and last decoder layer are shown in (b) and (c).

**Table 4.** Spatial and temporal dimensions. Keeping spatial-temporal feature dimensions and performing instance feature capture independently on different frames brings about 7.4 AP gains. Multi-scale feature maps can further bring 2.6 AP.

| Feature | AP | $AP_{50}$ | $AP_{75}$ | $AR_1$ | $AR_{10}$ |
|---|---|---|---|---|---|
| flatten | 35.1 | 56.8 | 35.6 | 38.1 | 41.8 |
| single-scale | 42.5 | 64.6 | 46.5 | 41.5 | 50.9 |
| multi-scale | 45.1 | 66.9 | 50.5 | 45.6 | 54.6 |

**Aggregation of temporal information.** The frame-level box queries and the predicted boxes can align the instance features from all frames, there are several ways to aggregate the aligned features into the instance query. We conduct an experiment to evaluate the different aggregation ways for these features, as shown in Table 5. In the 'sum' setting, the features from different frames are directly added together as the instance feature of this decoder layer. In the 'average' setting, the feature on each frame is averaged as the instance feature. In the 'weighted-sum' setting, we apply a softmax layer and a fully-connected layer on box embeddings to get the weights of each frame, and the features are aggregated in a weighted sum in Eq. 3. The result is 30.6 AP and 43.2 AP for 'sum' and 'average' settings respectively. Direct summation will cause the value to be unstable with different frame numbers. Since some instances only appear in a few frames, directly averaging features from all frames may cause the information to be diluted. Please refer to the Sup. Mat. for more details and visualization of different frame weights.

**Table 5.** Temporal information aggregation. Weighted sum brings a performance gain of 1.9 in AP.

| Aggregation | AP | $AP_{50}$ | $AP_{75}$ | $AR_1$ | $AR_{10}$ |
|---|---|---|---|---|---|
| sum | 30.6 | 44.5 | 34.3 | 37.2 | 45.0 |
| average | 43.2 | 65.2 | 48.5 | 43.4 | 52.8 |
| weighted-sum | 45.1 | 66.9 | 50.5 | 45.6 | 54.6 |

**Table 6.** Fewer frames for instance representation.We evenly sample fewer frames from a video to generate the mask head.

| Frames | AP | $AP_{50}$ | $AP_{75}$ | $AR_1$ | $AR_{10}$ |
|---|---|---|---|---|---|
| 1 | 38.1 | 58.3 | 41.3 | 38.7 | 47.5 |
| 3 | 43.4 | 65.4 | 47.6 | 42.4 | 51.3 |
| 5 | 44.6 | 66.5 | 49.7 | 44.8 | 54.6 |
| 10 | 44.7 | 66.9 | 49.5 | 44.3 | 53.5 |
| all | 45.1 | 66.9 | 50.5 | 45.6 | 54.6 |

**Robust instance representation.** Our decoder explicitly aligns and aggregates the information from each frame to learn a video-level instance representation. In this experiment, we try to generate instance representation with fewer frames. We use the instance representation to generate a mask head and apply the mask head on each frame to get the mask sequences, as shown in Table 6. Surprisingly, with only one frame as input, the generated mask head can produce a competitive result of 38.1 AP. With five frames as input, the performance is only 0.5 AP worse than taking all frames as input. This result shows that the mask head learned by our method can generalize well to unseen frames.

## 5    Conclusion

In this paper, we have proposed an effective transformer architecture for video instance segmentation, named SeqFormer, which performs attention mechanisms on each frame independently and learns a shared powerful instance query for each video-level instance. With the proposed instance query decomposition, our network can align the instance features and naturally tackle the instance tracking without additional tracking branches or post-processing. We demonstrated that our method surpasses all state-of-the-art methods by a large margin. We believe that our neat and efficient approach will serve as a strong baseline for future research in video instance segmentation.

# References

1. Arnab, A., Dehghani, M., Heigold, G., Sun, C., Lučić, M., Schmid, C.: Vivit: A video vision transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6836–6846 (2021) 4
2. Athar, A., Mahadevan, S., Osep, A., Leal-Taixé, L., Leibe, B.: Stem-seg: Spatiotemporal embeddings for instance segmentation in videos. In: ECCV (2020) 2, 4, 9, 10
3. Bertasius, G., Torresani, L.: Classifying, segmenting, and tracking object instances in video with mask propagation. In: CVPR (2020) 2, 4, 9, 10
4. Cao, J., Anwer, R.M., Cholakkal, H., Khan, F.S., Pang, Y., Shao, L.: Sipmask: Spatial information preservation for fast image and video instance segmentation. In: ECCV (2020) 1, 4, 10, 11
5. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV (2020) 2, 4, 5, 7, 8
6. Chen, X., Girshick, R., He, K., Dollár, P.: Tensormask: A foundation for dense object segmentation. In: ICCV (2019) 3
7. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020) 4
8. Fang, Y., Yang, S., Wang, X., Li, Y., Fang, C., Shan, Y., Feng, B., Liu, W.: Instances as queries. In: ICCV (2021) 3
9. Fu, Y., Yang, L., Liu, D., Huang, T.S., Shi, H.: Compfeat: Comprehensive feature aggregation for video instance segmentation. arXiv preprint arXiv:2012.03400 (2020) 10
10. Goel, V., Li, J., Garg, S., Maheshwari, H., Shi, H.: Msn: Efficient online mask selection network for video instance segmentation. arXiv preprint arXiv:2106.10452 (2021) 1
11. He, K., Gkioxari, G., Dollar, P., Girshick, R.: Mask R-CNN. In: ICCV (2017) 1, 3, 4
12. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) 5, 9
13. Huang, Z., Huang, L., Gong, Y., Huang, C., Wang, X.: Mask scoring r-cnn. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6409–6418 (2019) 3
14. Hwang, S., Heo, M., Oh, S.W., Kim, S.J.: Video instance segmentation using interframe communication transformers. In: NeurIPS (2021) 2, 4, 10, 11, 12
15. Kuhn, H.W.: The hungarian method for the assignment problem. Naval research logistics quarterly **2**(1-2), 83–97 (1955) 8
16. Li, K., Wang, Y., Gao, P., Song, G., Liu, Y., Li, H., Qiao, Y.: Uniformer: Unified transformer for efficient spatiotemporal representation learning. arXiv preprint arXiv:2201.04676 (2022) 4
17. Li, M., Li, S., Li, L., Zhang, L.: Spatial feature calibration and temporal fusion for effective one-stage video instance segmentation. In: CVPR (2021) 4
18. Lin, H., Wu, R., Liu, S., Lu, J., Jia, J.: Video instance segmentation with a propose-reduce paradigm. arXiv preprint arXiv:2103.13746 (2021) 4, 9, 10
19. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV (2017) 8

20. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014) 9
21. Liu, D., Cui, Y., Tan, W., Chen, Y.: Sg-net: Spatial granularity network for one-stage video instance segmentation. In: CVPR (2021) 1, 10
22. Liu, S., Qi, L., Qin, H., Shi, J., Jia, J.: Path aggregation network for instance segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8759–8768 (2018) 3
23. Liu, X., Wang, Q., Hu, Y., Tang, X., Zhang, S., Bai, S., Bai, X.: End-to-end temporal action detection with transformer. IEEE Transactions on Image Processing (TIP) (2022) 4
24. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. arXiv preprint arXiv:2103.14030 (2021) 10
25. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: ICLR (2019) 9
26. Meinhardt, T., Kirillov, A., Leal-Taixe, L., Feichtenhofer, C.: TrackFormer: Multi-object tracking with transformers. arXiv preprint arXiv:2101.02702 (2021) 4
27. Milletari, F., Navab, N., Ahmadi, S.A.: V-net: Fully convolutional neural networks for volumetric medical image segmentation. In: 2016 fourth international conference on 3D vision (3DV) (2016) 8
28. Nguyen, T.C., Tang, T.N., Phan, N.L., Nguyen, C.H., Yamazaki, M., Yamanaka, M.: 1st place solution for youtubevos challenge 2021: Video instance segmentation. arXiv preprint arXiv:2106.06649 (2021) 1
29. Patrick, M., Campbell, D., Asano, Y.M., Metze, I.M.F., Feichtenhofer, C., Vedaldi, A., Henriques, J., et al.: Keeping your eye on the ball: Trajectory attention in video transformers. arXiv preprint arXiv:2106.05392 (2021) 2, 4
30. Qi, J., Gao, Y., Hu, Y., Wang, X., Liu, X., Bai, X., Belongie, S., Yuille, A., Torr, P.H., Bai, S.: Occluded video instance segmentation: A benchmark. International Journal of Computer Vision pp. 1–18 (2022) 1
31. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: NeurIPS (2015) 3
32. Rezatofighi, H., Tsoi, N., Gwak, J., Sadeghian, A., Reid, I., Savarese, S.: Generalized intersection over union: A metric and a loss for bounding box regression. In: CVPR (2019) 8
33. Sun, P., Cao, J., Jiang, Y., Zhang, R., Xie, E., Yuan, Z., Wang, C., Luo, P.: Transtrack: Multiple object tracking with transformer. arXiv preprint arXiv:2012.15460 (2020) 4
34. Sun, P., Zhang, R., Jiang, Y., Kong, T., Xu, C., Zhan, W., Tomizuka, M., Li, L., Yuan, Z., Wang, C., et al.: Sparse r-cnn: End-to-end object detection with learnable proposals. In: CVPR (2021) 3, 4
35. Tian, Z., Shen, C., Chen, H.: Conditional convolutions for instance segmentation. In: ECCV (2020) 3, 7, 8
36. Tian, Z., Shen, C., Chen, H., He, T.: FCOS: Fully convolutional one-stage object detection. In: ICCV (2019) 4
37. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: NeurIPS (2017) 2, 4
38. Wang, X., Kong, T., Shen, C., Jiang, Y., Li, L.: Solo: Segmenting objects by locations. In: ECCV (2020) 3
39. Wang, Y., Xu, Z., Wang, X., Shen, C., Cheng, B., Shen, H., Xia, H.: End-to-end video instance segmentation with transformers. In: CVPR (2021) 2, 4, 6, 8, 10, 12

40. Wu, J., Jiang, Y., Sun, P., Yuan, Z., Luo, P.: Language as queries for referring video object segmentation. In: CVPR. pp. 4974–4984 (2022) 4
41. Xie, E., Sun, P., Song, X., Wang, W., Liu, X., Liang, D., Shen, C., Luo, P.: Polarmask: Single shot instance segmentation with polar representation. In: CVPR. pp. 12193–12202 (2020) 3
42. Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers. Advances in Neural Information Processing Systems **34**, 12077–12090 (2021) 4
43. Xu, N., Yang, L., Yang, J., Yue, D., Fan, Y., Liang, Y., Huang., T.S.: Youtubevis dataset 2021 version. https://youtube-vos.org/dataset/vis/ 8
44. Yan, B., Peng, H., Fu, J., Wang, D., Lu, H.: Learning spatio-temporal transformer for visual tracking. In: ICCV (2021) 4
45. Yang, L., Fan, Y., Xu, N.: Video instance segmentation. In: ICCV (2019) 1, 3, 4, 8, 9, 10, 11
46. Yang, S., Fang, Y., Wang, X., Li, Y., Fang, C., Shan, Y., Feng, B., Liu, W.: Crossover learning for fast online video instance segmentation. arXiv preprint arXiv:2104.05970 (2021) 1, 4, 10, 11
47. Zhao, Y., Xiong, Y., Lin, D.: Trajectory convolution for action recognition. In: NeurIPS (2018) 2, 12
48. Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng, J., Xiang, T., Torr, P.H., 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. pp. 6881–6890 (2021) 4
49. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable detr: Deformable transformers for end-to-end object detection. arXiv preprint arXiv:2010.04159 (2020) 3, 4, 5, 6, 9
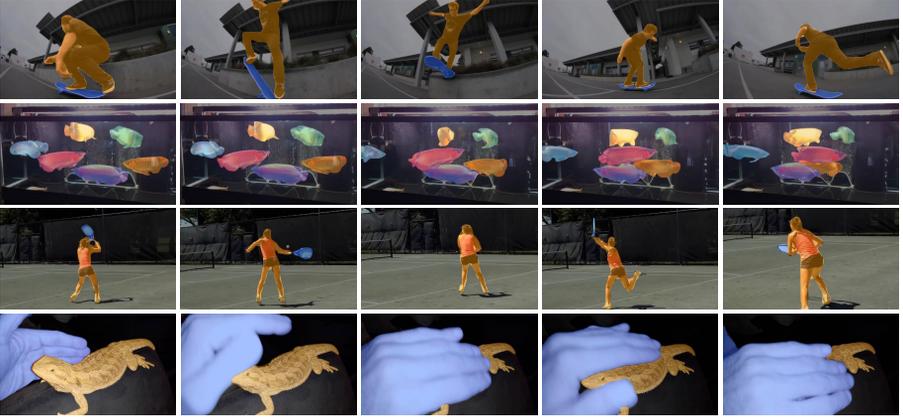
# A    Appendix



**Fig. 5.** Visualization of SeqFormer on the YouTube-VIS 2019 validation dataset. The first row shows the instances with various poses. The second row shows the case of a lot of similar instances that are close together with overlapping. The third row shows the situation where an instance reappears after being occluded while in motion. The last row shows an instance severely occluded by the other instance. The same colors depict the mask sequences of the same instances

## A.1    Visualization

In Fig. 5, we visualize the results of SeqFormer with four challenging cases. It can be seen that SeqFormer can handle these situations well. In Fig. 6, we show more qualitative results of the intermediate attention of transformer decoder. Since the same initial instance query is used to predict sampling points for each frame in the first decoder layer (Eq.1), the distribution of sampling points on each frame is the same in Fig. 6 (a) and (d). After that, the initial instance query is decomposed into frame-level box queries that are kept and maintained independently on each frame. Starting from the second layer of the SeqFormer decoder, the box query is used to predict the sampling points of the current frame, and the sampled features are used to refine the box query for the next decoder layer. By doing so, SeqFormer attends to different spatial locations following the motion of the instance in a coarse-to-fine manner.

## A.2    Aggregation of Temporal Information

SeqFormer is able to attend to different spatial locations following the motion of the instance. The aligned features are aggregated into an instance query
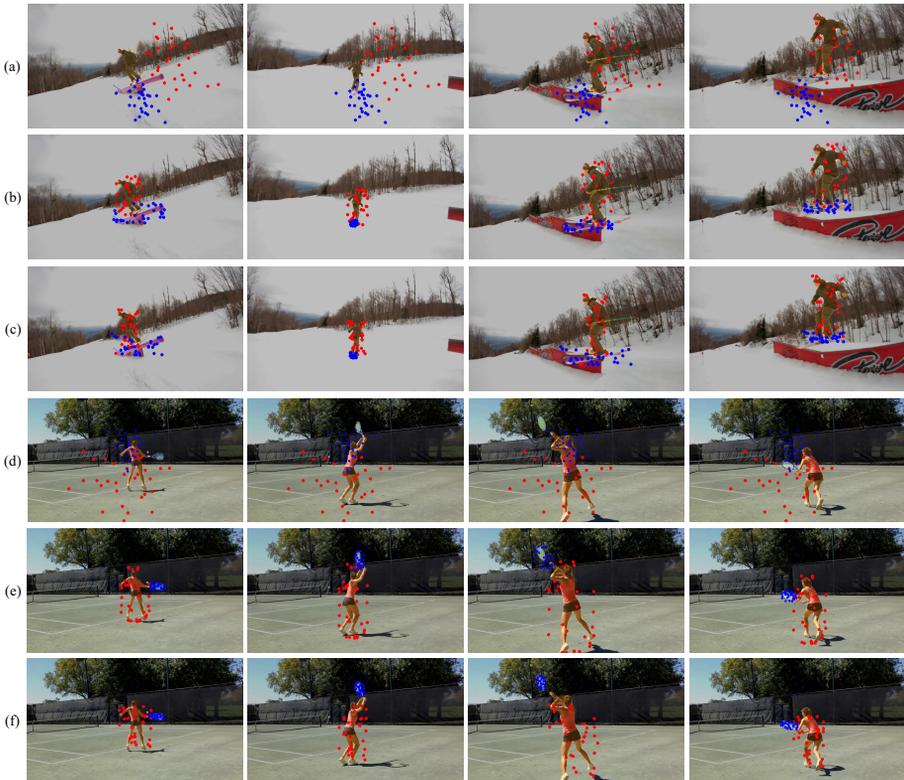
**Fig. 6.** Visualization of attention. We draw the sampling points that the deformable attention attends to. The four frames in each row are from the same video. Each sampling point is marked as a filled circle whose color indicates its corresponding instance identity. (a) and (d) show the sampling points from the first decoder layer. (b) and (e) show the sampling points from the second decoder layer. (c) and (f) show the sampling points from the last decoder layer.

to generate a video-level instance representation. However, an instance may not appear in every frame due to occlusion and camera motion. The features from frames without instance are useless or even harmful. To address this, SeqFormer aggregates temporal features in a weighted manner, where the weights are learned upon the box queries in Eq.3. We visualize the learned weights and the corresponding frames in Fig. 7. It can be seen that the features from frames without instance have lower weights.

## A.3   Qualitative Comparisons

We provide some qualitative comparisons with other methods in Fig. 8, the mask predictions of SeqFormer are more stable over time. More video results and comparisons can be found in the rest of the supplementary material.
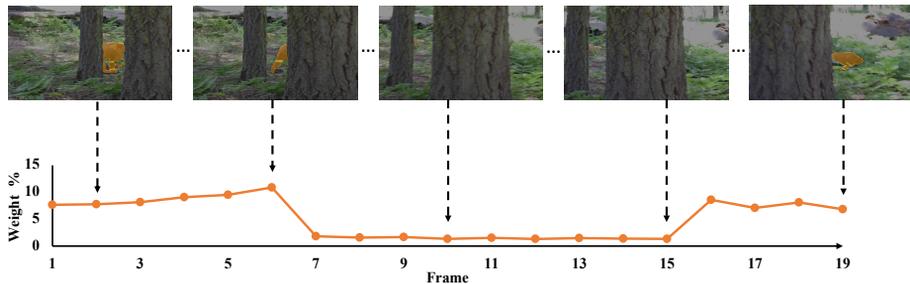
**Fig. 7.** Visualization of the normalized softmax weights and the corresponding frames.

## A.4   Clip Matching

Our model can be extended to per-clip model through clip matching algorithm to handle long videos. Specifically, we divide long videos into clips with overlapping frames, and match clip-level instance masks by calculate the matching scores which are space-time soft IoU of overlapping frames, following IFC. We evaluate our method by varying the length of clips in Table 7. Our method still achieves competitive performance but slightly worse when evaluating in the clip-wise manner. This manner makes our method handle very long videos with limited computational resources and have a wider range of application scenarios.

**Table 7.** Evaluating SeqFormer in a clip-wise manner.

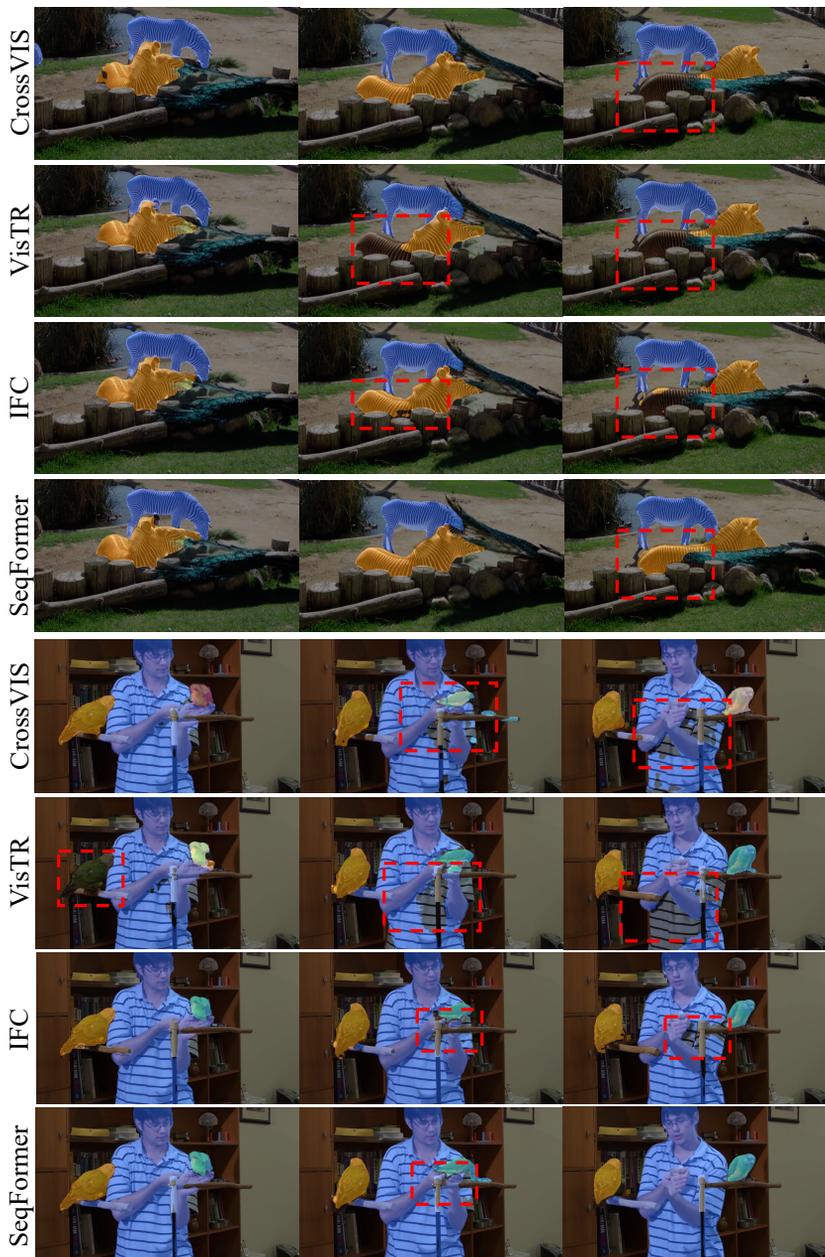| Clip Length | Whole | 20 | 15 | 10 |
|---|---|---|---|---|
| AP | 45.1 | 44.8 | 43.6 | 42.0 |

**Fig. 8.** Qualitative comparisons with other methods on YouTube-VIS 2019. All methods use ResNet-50 backbone. The three frames in each row are from the same video. The mask predictions of SeqFormer are more stable over time. Best viewed in color.