# Video Object Segmentation with Language Referring Expressions

Anna Khoreva[1,2], Anna Rohrbach[3], and Bernt Schiele[1]

[1]Max Planck Institute for Informatics    [2]Bosch Center for Artificial Intelligence
[3]University of California, Berkeley

**Abstract** Most state-of-the-art semi-supervised video object segmentation methods rely on a pixel-accurate mask of a target object provided for the first frame of a video. However, obtaining a detailed segmentation mask is expensive and time-consuming. In this work we explore an alternative way of identifying a target object, namely by employing language referring expressions. Besides being a more practical and natural way of pointing out a target object, using language specifications can help to avoid drift as well as make the system more robust to complex dynamics and appearance variations. Leveraging recent advances of language grounding models designed for images, we propose an approach to extend them to video data, ensuring temporally coherent predictions. To evaluate our approach we augment the popular video object segmentation benchmarks, DAVIS$_{16}$ and DAVIS$_{17}$ with language descriptions of target objects. We show that our language-supervised approach performs on par with the methods which have access to a pixel-level mask of the target object on DAVIS$_{16}$ and is competitive to methods using scribbles on the challenging DAVIS$_{17}$ dataset.

*Query: "A man in a red sweatshirt performing breakdance"*



Figure 1: Examples of the proposed approach. Classical semi-supervised video object segmentation relies on an expensive pixel-level mask annotation of a target object in the first frame of a video. We explore a more natural and more practical way of pointing out a target object by providing a language referring expression.

## 1 Introduction

Video object segmentation has recently witnessed growing interest [3,15,6,41]. Segmenting objects at pixel level provides a finer understanding of video and is relevant for many applications, e.g. augmented reality, video editing, and rotoscoping.

Ideally, one would like to obtain a pixel-accurate segmentation of objects in video with no human input during test time. However, the current state-of-the-art unsupervised video object segmentation methods [54,17,47] have troubles segmenting the target objects in videos containing multiple objects and cluttered backgrounds without any guidance from the user. Hence, many recent works [3,15,50] employ a semi-supervised approach, where a pixel-level mask of the target object is manually annotated in the first frame and the task is to accurately segment the object in successive frames. Although

this setting has proven to be successful, it can be prohibitive for many applications. It is tedious and time-consuming for the user to provide a pixel-accurate segmentation and usually takes more than a minute to annotate a single instance ([26] reports 79s for polygon annotations, precisely delineating an object would take even more). To make video object segmentation more applicable in practice, instead of costly pixel-level masks [41,32,2] propose to employ point clicks or scribbles to specify the target object in the first frame. This is much faster and takes an annotator on average 7.5s to label an object with point clicks [32] and 10s with scribbles [25]. However, on small touchscreen devices, such as tablets or phones, providing precise clicks or drawing scribbles using fingers could be cumbersome and inconvenient for the user.

To overcome these limitations we propose a new task - segmenting objects in video using language referring expressions - which is a more natural way of human-computer interaction. It is much easier for a user to say: "Segment the man in a red sweatshirt performing breakdance" (see Figure 1), than to provide a tedious pixel-level segmentation mask or struggle with drawing a scribble which does not straddle the object boundary. Moreover, employing language specifications can make the system more robust to background clutter, help to avoid drift and better adapt to the complex dynamics inherent to videos, while not over-fitting to a particular view in the first frame (see Table 4).

We aim to investigate the capabilities and limitations of existing techniques on the proposed task and explore how far one can go while leveraging the advances in image-level language grounding and pixel-level segmentation in videos. We start by analyzing the performance of the state-of-the-art language grounding models [58,56] for localization of objects in videos via bounding boxes. We discover that they suffer from a number of issues, predicting temporally inconsistent and jittery boxes, and show a way to enhance their predictions by enforcing temporal coherency (see Figure 3). Next we propose a convnet-based framework that utilizes referring expressions for video object segmentation task, where the output of the grounding model (bounding box) is used as a guidance for pixel-wise segmentation of the object. We also show that video object segmentation using the mask annotation on the first frame can be further improved by using language supervision, highlighting the complementarity of both modalities.

To evaluate the proposed approach we extend the popular benchmarks for segmenting single and multiple objects in videos, DAVIS$_{16}$ [38] and DAVIS$_{17}$ [42], with language descriptions of the target objects. We collect the annotations using two different settings, asking the annotators to provide a description of the target object based on the first frame only as well as on the full video. Future work may choose which setting they prefer to use. On average each video has been annotated with 7.5 referring expressions and it takes the annotator around 5s to provide a referring expression for a target object.

Our language-supervised approach performs on par with semi-supervised methods which have access to the pixel-accurate object mask on DAVIS$_{16}$ and shows comparable results to the techniques that employ scribbles on the challenging DAVIS$_{17}$ dataset.

In summary, our contributions are the following. We present a new task of segmenting objects in video using natural language referring expressions for which we augment two well-known video segmentation benchmarks with textual descriptions of target objects. We conduct an extensive analysis of the performance of the state-of-the-art language grounding models on video data and propose a way to improve their

temporal coherency. To the best of our knowledge we are the first to perform an analysis of transferability of image-based grounding models to video. We show that high quality video object segmentation results can be obtained by employing language referring expressions, allowing a more natural and practical human-computer interaction. Moreover, we show that language descriptions are complementary to visual forms of supervision, such as masks, and can be exploited as an additional source of guidance for object segmentation. Thus, while proposing the new task and accompanying dataset, our work contributes the necessary benchmark analysis, a very competitive baseline and valuable insights for future work. We hope our findings would further promote the research in the field of video object segmentation via language expressions and help to discover better techniques that can be used in realistic scenarios.

## 2    Related Work

### 2.1    Grounding natural language expressions

There has been an increasing interest in the task of grounding natural language expressions over the last few years [57,27,23]. We group the existing works by the type of visual domain: images and video.

**Image domain.** Grounding natural language expressions is a task of localizing a given expression in an image with a bounding box [58,34] or a segmentation mask [27,23]. Referring expression comprehension is a closely related task, where the goal is to localize the non-ambiguous referring expression. Most existing approaches rely on external bounding box proposals which are scored to determine the top scoring box as the correct region [30,56]. A few recent works explore methods of inferring object regions by proposal generation network [4] or efficient subwindow search [55]. Multiple existing approaches model relationships between objects present in the scene [35,14]. In this work we choose two state-of-the-art grounding models for experimentation and analysis [58,56]. DBNet [58] frames grounding as a classification task, where an expression and an image region serve as input and a binary classification decision is an output. A key component of this approach is utilization of negative expressions and image regions to ensure discriminative training. DBNet currently leads on Visual Genome [22]. MattNet [56] is a modular network which "softly" decomposes referring expressions in three parts: subject, location, and relationship, each of which is processed by a different visual module. This allows MattNet to process referring expressions of general forms, as each module can be "enabled" or "disabled" depending on the expression. MattNet achieves top performance on RefCOCO(g/+) [57,34] both in terms of bounding box localization and pixel-wise segmentation accuracy.

**Video domain.** The progress made in image-level natural language grounding leads to an increasing interest in application to video. The recent work of [24] studies object tracking in video using language expressions. They introduce a dynamic convolutional layer, where a language query is used to predict visual convolutional filters. [1] addresses object tracking in video with the language descriptions and human gaze as input. Our work falls in the same line of research, as we are exploring natural language as input for video object segmentation. To the best of our knowledge, this is the first work to apply natural language to this task. A concurrent work by [10] has addressed a task of actor/action segmentation in video based on sentence input. Their work focuses

on seven classes of actors (adult, baby, etc.) and mostly action-oriented descriptions. In contrast, we consider arbitrary objects and unconstrained referring expressions.

### 2.2   Video Object Segmentation

Video object segmentation has witnessed considerable progress [37,48,47,21,3,50]. In the following, we group the related work into unsupervised and semi-supervised.

**Unsupervised methods.**  Unsupervised methods assume no human input on the video during test time. They aim to group pixels consistent in both appearance and motion and extract the most salient spatio-temporal object tube. Several techniques exploit object proposals [54,21], saliency [9] and optical flow [37]. Convnet-based approaches [6,17,47] cast video object segmentation as a foreground/background classification problem and feed to the network both appearance and motion cues. Because these methods do not have any knowledge of the target object, they have difficulties in videos with multiple moving and dominant objects and cluttered backgrounds.

**Semi-supervised methods.**  Semi-supervised methods assume human input for the first frame, either by providing a pixel-accurate mask [48,3], clicks [32] or scribbles [41], and then propagate the information to the successive frames. Existing approaches focus on leveraging superpixels [53], constructing graphical models [48], utilizing object proposals [40] or employing optical flow and long-term trajectories [52]. Lately, convnets have been considered for the task [3,39,50]. These methods usually build the architecture upon the semantic segmentation networks [29] and process each frame of the video individually. [3] proposes to fine-tune a pre-trained generic object segmentation network on the first frame mask of the test video to make it sensitive to the target object. [39] employs a similar strategy, but also provides a temporal context by feeding the previous frame mask to the network. Several methods extend the work of [3] by incorporating the semantic information [33] or by integrating online adaptation [50]. [15] proposes to employ a recurrent network to exploit the long-term temporal information.

The above methods employ a pixel-level mask on the first frame. However, for many applications, particularly on small touchscreen devices, it can be prohibitive to provide a pixel-accurate segmentation. Hence, there has been a growing interest to integrate cheaper forms of supervision, such as point clicks [2,32] or scribbles [41], into convnet-based techniques. In spirit with these approaches, we aim to reduce the annotation effort on the first frame by using language referring expressions to specify the object. Our approach also builds upon convnets and exploits both linguistic and visual modalities.

## 3   Method

In this section we provide an overview of the proposed approach. Given a video $V = \{f_1, ..., f_N\}$ with N frames and a textual query of the target object $Q$, our aim is to obtain a pixel-level segmentation mask of the target object in every frame that it appears.

We leverage recent advances in grounding referring expressions in images [58,56] and pixel-level segmentation in videos [39,17]. Our method consists of two main steps (see Figure 2). Using as input the textual query $Q$ provided by the user, we first generate target object bounding box proposals for every frame of the video by exploiting referring expression grounding models, designed for images only. Applying these models off-the-shelf results in temporally inconsistent and jittery box predictions (see Figure
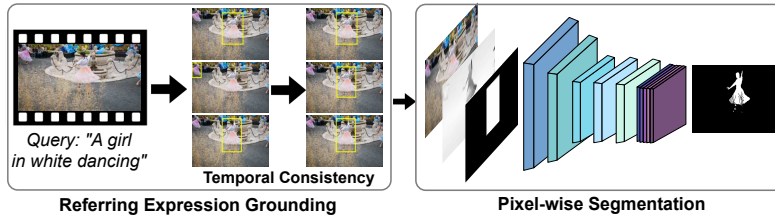
Figure 2: System overview. We first localize the target object via grounding model using the given referring expression and enforce temporal consistency of bounding boxes across frames. Next we apply a segmentation convnet to recover detailed object masks.

3). Therefore, to mitigate this issue and make them more applicable for video data, we next employ temporal consistency, which enforces bounding boxes to be coherent across frames. As a second step, using as guidance the obtained box predictions of the target object on every frame of the video we apply a convnet-based pixel-wise segmentation model to recover detailed object masks in each frame.

### 3.1   Grounding objects in video by referring expressions

As discussed in §2, the task of natural language grounding is to automatically localize a region described by a given language expression. It is typically formulated as measuring the compatibility between a set of object proposals $O = \{o_i\}_{i=1}^{M}$ and a given textual query $Q$. The grounding model provides as output a set of matching scores $S = \{s_i\}_{i=1}^{M}$ between a box proposal and a textual query $Q$. The box proposal with the highest matching score is selected as the predicted region.

We employ two state-of-the-art referring expression grounding models – DBNet [58] and MattNet [56], to localize the object in each frame. Mask R-CNN [12] bounding box proposals are exploited as an initial set of proposals for both models, although originally DBNet has been designed to utilize EdgeBox proposals [8]. However, using the grounding models designed for images and picking the highest scoring proposal for each video frame lead to temporally incoherent results. Even with simple textual queries for adjacent frames that from a human perspective look very much alike, the referring model often outputs inconsistent predictions (see Figure 3). This indicates the inherent instability of the grounding models trained on the image domain. To resolve this problem we propose to re-rank the object proposals by exploiting temporal structure along with the original matching scores given by a grounding model.

**Temporal consistency.**  The goal of the temporal smoothing step is to improve temporal consistency and to reduce id-switches for target object predictions across frames. Since objects tend to move smoothly through space and in time, there should be little changes from frame to frame and the box proposals should have high overlap between neighboring frames. By finding temporally coherent tracks of an object that are spread-out in time, we can focus on the predictions that consistently appear throughout the video and give less emphasis to objects that appear for only a short period of time.

The grounding model provides the likeliness of each box proposal to be the target object by outputting a matching score $s_i$. Then each box proposal is re-ranked based on its overlap with the proposals in other frames, the original objectness score given by [12] and its matching score from the grounding model. Specifically, for each proposal

we compute a new score: $\hat{s}_i = s_i * (\sum_{j=1, j \neq i}^{M} r_{ij} * d_j * s_j / t_{ij})$, where $r_{ij}$ measures an intersection-over-union ratio between box proposals $i$ and $j$, $t_{ij}$ denotes the temporal distance between two proposals ($t_{ij} = |f_i - f_j|$) and $d_j$ is the original objectness score. Then, in each frame we select the proposals with the highest new score. The new scoring rewards temporally coherent predictions which likely belong to the target object and form a spatio-temporal tube. This step allows to improve temporal coherence boosting grounding and video segmentation performance (see Table 1 in §5 and Table 5 in §6) while being computational efficient (takes only a fraction of second).

### 3.2   Pixel-level video object segmentation

We next show how to output pixel-level object masks, exploiting the bounding boxes from grounding as a guidance for the segmentation network. The boxes are used as the input to the network to guide the network towards the target object, providing its rough location and extent. The task of the network is to obtain a pixel-level foreground/background segmentation mask using appearance and motion cues.

**Approach.** We model pixel-level segmentation as a box refinement task. The bounding box is transformed into a binary image (255 for the interior of the box, 0 for the background) and concatenated with the RGB channels of the input image and optical flow magnitude, forming a 5-channel input for the network. Thus we ask the network to learn to refine the provided boxes into accurate masks. Fusing appearance and motion cues allows to better exploit video data and handle better both static and moving objects.

We make one single pass over the video, applying the model per-frame. The network does not keep a notion of the specific appearance of the object in contrast to [39,3], where the model is fine-tuned during the test time to learn the appearance of the target object. Neither do we do an online adaptation as in [50], where the model is updated on its previous predictions while processing video frames. This makes the system more efficient during the inference time, which is more suitable for real-world applications.

Similar to [39], we train the network on static images, employing the saliency segmentation dataset [7] which contains a diverse set of objects. The bounding box is obtained from the ground truth masks. To make the system robust during test time to sloppy boxes from the grounding model, we augment the ground truth box by randomly jittering its coordinates (uniformly, $\pm 20\%$ of the original box width and height). We synthesize optical flow from static images by applying affine transformations for both background and foreground object to simulate the camera and object motion in the neighboring frames, as in [20]. This simple strategy allows us to train on diverse set of static images, while exploiting motion information during test time. We train the network on many triplets of RGB images, synthesized flow magnitude images and loose boxes in order for the model generalize well to different localization quality of boxes given by the grounding model and different dynamics of the object.

During inference we use the state-of-the-art optical flow estimation method Flow-Net2.0 [16]. We compute the optical flow magnitude by subtracting the median motion for each frame and averaging the magnitude of the forward and backward flow. The obtained image is further scaled to [0; 255] to maintain the same range as RGB channels.

**Network.** As our network architecture we use ResNet-101 [13]. We adapt the network to the segmentation task following the procedure of [29] and employing atrous convolutions [5] with hybrid rates [51] within the last two blocks of ResNet to enlarge the
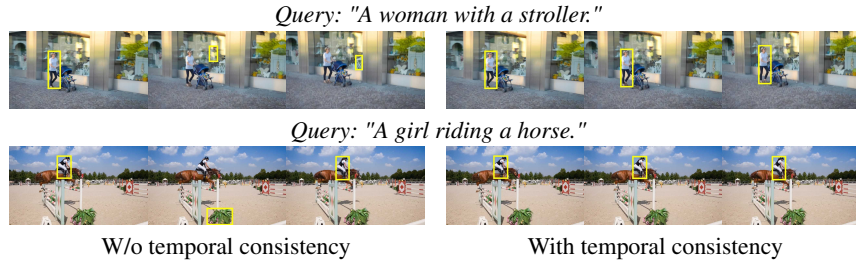
*Query: "A woman with a stroller."*



*Query: "A girl riding a horse."*



W/o temporal consistency                    With temporal consistency

Figure 3: Qualitative results of language grounding with and w/o temporal consistency on DAVIS$_{17}$. The results are obtained using MattNet [56] trained on RefCOCO [57].

receptive field as well as to alleviate the "gridding" issue. After the last block, we apply spatial pyramid pooling [5], which aggregates features at multiple scales by applying atrous convolutions with different rates, and augment it with the image-level features [28] to exploit better global context. The network is trained using a standard cross-entropy loss (all pixels are equally weighted). The final logits are upsampled to the ground truth resolution to preserve finer details for back-propagation.

For network initialization we use a model pre-trained on ImageNet [13]. The new layers are initialized using the "Xavier" strategy [11]. The network is trained on MSRA [7] for segmentation. To avoid the domain shift we fine-tune the model on the training sets of DAVIS$_{16}$ [38] and DAVIS$_{17}$ [42] respectively. We employ SGD with a polynomial learning policy with initial learning rate of $0.001$, crop size of $513 \times 513$, random scale data augmentation (from $0.5$ to $2.0$) and left-right flipping during training. The network is trained for $20k$ iterations on MSRA and $20k$ iterations on the training set of DAVIS$_{16}$/DAVIS$_{17}$. During inference we employ test time augmentation as in [5].

**Other sources of supervision.** Additionally we consider variants of the proposed model using different sources of supervision. Our approach is flexible and can take advantage of the first frame mask annotation as well as language. We describe how language can be used on top of the mask supervision, improving the robustness of the system against occlusions and dynamic backgrounds (see §6 for results).

*Mask.* Here we discuss a variant that uses only the first frame mask supervision during test time. The network is initialized with the bounding box obtained from the object mask in the 1st frame and for successive frames uses the prediction from the preceding frame warped with the optical flow (as in [39]) to get the input box for the next frame. Following [39,3] we fine-tune the model for $1k$ iterations on an augmented set obtained from the first frame image and mask, to learn the specific properties of the object.

*Mask + Language.* We show that using language supervision is complementary to the first frame mask. Instead of relying on the preceding frame prediction as in the previous paragraph, we use the bounding boxes obtained from the grounding model after the temporal consistency step. We initialize with the ground truth box in the first frame and fine-tune the network on the 1st frame.

## 4    Collecting referring expressions for video

Our task is to localize and provide a pixel-level mask of an object on all video frames given a language referring expression obtained either by looking at the first frame only

*ID 1: "A man in a grey t-shirt and yellow trousers"*
*ID 2: "A woman in a black shirt"*
*ID 3: "A white truck on the road"*
First frame annotation

*ID 1: "A man in a grey shirt walking through the crossing"*
*ID 2: "A woman walking through the crossing"*
*ID 3: "A white truck moving from the left to right"*
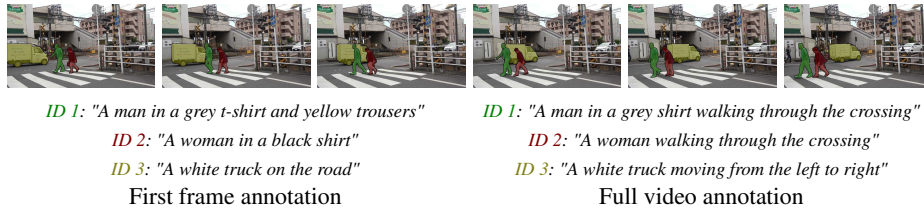Full video annotation

Figure 4: Example of annotations provided for the 1st frame vs. the full video. Full video annotations include descriptions of activities and overall are more complex.

or the full video. To validate our approach we employ two popular video object segmentation datasets, DAVIS$_{16}$ [38] and DAVIS$_{17}$ [42]. These two datasets introduce various challenges, containing videos with single or multiple salient objects, crowded scenes, similar looking instances, occlusions, camera view changes, fast motion, etc.

DAVIS$_{16}$ [38] consists of 30 training and 20 test videos of diverse object categories with all frames annotated with pixel-level accuracy. Note that in this dataset only a single object is annotated per video. For the multiple object video segmentation task we consider DAVIS$_{17}$. Compared to DAVIS$_{16}$, this is a more challenging dataset, with multiple objects annotated per video and more complex scenes with more distractors, occlusions, smaller objects, and fine structures. Overall, DAVIS$_{17}$ consists of a training set with 60 videos, and a validation/test-dev/test-challenge set with 30 sequences each.

As our goal is to segment objects in videos using language specifications, we augment all objects annotated with mask labels in DAVIS$_{16}$ and DAVIS$_{17}$ with non-ambiguous referring expressions. We follow the work of [34] and ask the annotator to provide a language description of the object, which has a mask annotation, by looking only at the first frame of the video. Then another annotator is given the first frame and the corresponding description, and asked to identify the referred object. If the annotator is unable to correctly identify the object, the description is corrected to remove ambiguity and to specify the object uniquely. We have collected two referring expressions per target object annotated by non-computer vision experts (Annotator 1, 2).

However, by looking only at the 1st frame, the obtained referring expressions may potentially be invalid for an entire video. (We actually quantified that only∼ 15% of the collected descriptions become invalid over time and it does not affect strongly segmentation results as temporal consistency step helps to disambiguate some of such cases, see the supp. material for details.) Besides, in many applications, such as video editing or video-based advertisement, the user has access to a full video. Providing a language query which is valid for all frames might decrease the editing time and result in more coherent predictions. Thus, on DAVIS$_{17}$ we asked the workers to provide a description of the object by looking at the full video. We have collected one expression of the full video type per target object. Future work may choose to use either setting.

The average length for the first frame/full video expressions is 5.5/6.3 words. For DAVIS$_{17}$ first frame annotations we notice that descriptions given by Annotator 1 are longer than the ones by Annotator 2 (6.4 vs. 4.6 words). We evaluate the effect of description length on the grounding performance in §5. Besides, the expressions relevant to a full video mention verbs more often than the first frame descriptions (44% vs. 25%). This is intuitive, as referring to an object which changes its appearance and

*ID 1: "A girl with blonde hair dressed in blue".*



*ID 1: "A brown camel in the front".*



*ID 1: "A black scooter ridden by a man". ID 2: "A man in a suit riding a scooter".*



Figure 5: Video object segmentation qualitative results using only referring expressions as supervision on DAVIS$_{16}$ and DAVIS$_{17}$, val sets. Frames sampled along the video.

position over time may require mentioning its actions. Adjectives are present in over $50\%$ for all annotations. Most of them refer to colors (over $70\%$), shapes and sizes ($7\%$) and spatial/ordering words ($6\%$ first frame vs. $13\%$ full video expressions). The full video expressions also have a higher number of adverbs and prepositions, and overall are more complex than the ones provided for the first frame, see Figure 4 for examples.

Overall augmented DAVIS$_{16/17}$ contains $\sim 1.2$k referring expressions for more than 400 objects on 150 videos with $\sim 10$k frames. We believe the collected data will be of interest to segmentation as well as vision and language communities, providing an opportunity to explore language as alternative input for video object segmentation.

## 5   Evaluation of natural language grounding in video

In this section we discuss the performance of natural language grounding models on video data. We experiment with DBNet [58] and MattNet [57]. DBNet is trained on Visual Genome [22] which contains images from MS COCO [26] and YFCC100M [45], and spans thousands of object categories. MattNet is trained on referring expressions for MS COCO images [26], specifically RefCOCO and RefCOCO+ [57]. Unlike RefCOCO which has no restrictions on the expressions, RefCOCO+ contains no spatial words and rather focuses on object appearance. Both aforementioned models rely on external bounding box proposals, such as EdgeBox [8] or Mask R-CNN [12].

We carry out most of our evaluation on DAVIS$_{16}$ and DAVIS$_{17}$ with the referring expressions introduced in §4. To evaluate the localization quality we employ the intersection-over-union overlap (IoU) of the top scored box proposal with the ground truth bounding box, averaged across all queries.

### 5.1   DAVIS$_{16}$/DAVIS$_{17}$ referring expression grounding

Table 1 reports performance of the grounding models on DAVIS$_{16}$ and DAVIS$_{17}$ referring expressions. In the following we summarize our key observations.

(1) We see the effect of replacing EdgeBox with Mask R-CNN object proposals for DBNet model ($54.1$ to $64.9$). Employing better proposals significantly improves the quality of this grounding method, thus we rely on Mask R-CNN proposals in all the following experiments. (2) We note the stability of grounding performance across two annotations (see $\Delta$(A1,A2)), showing that the grounding methods are quite robust to

| Method | Object proposals | Train. data | Temp. cons. | DAVIS$_{16}$ 1st frame | | DAVIS$_{17}$ 1st frame | | Full video |
| | | | | mIoU | $\Delta$(A1,A2) | mIoU | $\Delta$(A1,A2) | mIoU |
|---|---|---|---|---|---|---|---|---|
| DBNet | EdgeBox | Vis.Gen. | - | 54.1 | 1.0 | - | - | - |
| | Mask R-CNN | | - | 64.9 | 2.1 | 48.4 | 1.3 | 49.6 |
| MattNet | Mask R-CNN | RefCOCO | - | 67.1 | 2.2 | 51.6 | 1.6 | 50.3 |
| | | RefCOCO+ | - | 69.1 | 3.2 | 50.8 | 1.2 | 50.1 |
| DBNet | Mask R-CNN | Vis.Gen. | ✓ | 68.8 | 0.6 | 49.6 | 1.6 | 50.2 |
| MattNet | Mask R-CNN | RefCOCO | ✓ | 71.4 | 0.2 | 52.8 | 0.5 | 51.3 |
| | | RefCOCO+ | ✓ | 72.5 | 0.3 | 52.3 | 0.0 | 51.2 |

Table 1: Comparison of the DBNet[58] and MattNet [56] models on DAVIS$_{16}$ training set and DAVIS$_{17}$ val set. $\Delta$(A1,A2) denotes the difference between Annotator 1 and 2.

variations in language expressions. (3) The grounding models trained on images are not stable across frames, even when small changes in appearance occur (e.g. see Figure 3). We see that our proposed temporal consistency technique benefits both methods (e.g. DBNet: 64.9 vs. 68.8 on DAVIS$_{16}$, MattNet 51.6 vs. 52.8 on DAVIS$_{17}$). (4) On both datasets MattNet performs better than DBNet. The gap is particularly large on DAVIS$_{16}$ (72.5 vs. 68.8), as DAVIS$_{16}$ contains videos of a single foreground moving object, while DBNet is trained on a densely labeled Visual Genome dataset with many foreground and background objects. (5) On DAVIS$_{16}$ MattNet trained on RefCOCO+ outperforms MattNet trained on RefCOCO (72.5 vs. 71.4), while both perform similar on DAVIS$_{17}$. As RefCOCO+ contains no spatial words, MattNet trained on this dataset is more accurate in localizing queries mentioning object appearance. (6) Compared to DAVIS$_{16}$, DAVIS$_{17}$ is significantly more challenging, as it contains cluttered scenes with multiple moving objects (e.g. for MattNet 71.4 vs. 52.8). (7) When comparing results on expressions provided for the first frame versus expressions provided for the full video, we observe diverging trends. While DBNet is able to improve its performance (48.4 vs. 49.6), MattNet performance decreases (52.8 vs. 51.3). We attribute this to the fact that DBNet is trained on the more diverse Visual Genome descriptions.

**Attribute-based analysis.** Next we perform a more detailed analysis of the grounding models on DAVIS$_{17}$. We split the textual queries/videos into subsets where a certain attribute is present and report the averaged results for the subsets. Table 2 presents attribute-based grounding performance on first-frame based expressions averaged across annotators. To estimate the upper bound performance and the impact of imperfect bounding box proposals we add an Oracle comparison, where performance is reported on the ground-truth object boxes. We summarize our findings in the following.

(1) As MattNet is trained on MS COCO images and both models rely on MS COCO-based Mask R-CNN proposals, we compare performance for expressions which include COCO versus non-COCO objects. Both models drop in performance on non-COCO expressions, showing the impact of the domain shift to DAVIS$_{17}$ (e.g. for MattNet 59.6 vs. 36.9). Even DBNet which is trained on a larger training corpus suffers from the same effect (55.5 vs. 37.3). (2) We label the DAVIS$_{17}$ expressions as "spatial" if they include some of the spatial words (e.g. left, right). Such queries are significantly harder for all models (e.g. for MattNet 33.8 vs. 58.5). (3) Verbs are important as they allow to disambiguate an object in a video based on its actions. Presence of verbs in expres-

| Method | Train. data | Obj. prop. | mIoU | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | CO. | ˜CO. | Sp. | ˜Sp. | Ve. | ˜Ve. | \multicolumn{3}{c}{Expr. length} | | | \multicolumn{3}{c}{Num. obj.} |
| | | | | | | | | | S | M | L | 1 | 2-3 | >3 |
| DBNet | Vis.Gen. | Mask | 55.5 | 37.3 | **36.5** | 55.7 | 37.4 | **52.0** | 61.8 | 49.2 | 33.6 | 79.5 | 49.3 | **22.6** |
| MattNet | RefCOCO | R-CNN | **59.6** | 36.9 | 33.8 | **58.5** | **55.8** | 51.7 | **63.9** | **50.2** | **49.1** | **86.1** | **51.2** | 16.1 |
| DBNet | Vis.Gen. | Oracle | **79.3** | **59.0** | **47.7** | **81.7** | 70.3 | **77.6** | **84.8** | **69.9** | **67.9** | 100 | **73.8** | **37.2** |
| MattNet | RefCOCO | | 73.2 | 46.6 | 42.2 | 72.5 | **74.7** | 62.9 | 79.0 | 61.1 | 59.0 | 100 | 64.5 | 23.2 |

Table 2: Grounding performance breakdown for different attributes on DAVIS$_{17}$, val set. Results obtained after the temporal consistency, using average between two annotators (1st frame based). Attributes: COCO/non-COCO, Spatial/non-Spatial, Verbs/no Verbs, Expression length (Short, Medium, Long) and Number of objects.

sions is a challenging factor for DBNet trained on Visual Genome, while MattNet does significantly better (37.4 vs. 55.8). (4) Expression length is also an important factor. We quantize our expressions into Short (<4 words), Medium (4–6 words) and Long (>6 words). All models demonstrate similar drop in performance as expression length increases (e.g. for MattNet $63.9 \rightarrow 50.2 \rightarrow 49.1$). (5) Videos with more objects are more difficult, as these objects also tend to be very similar, such as e.g. fish in a tank (e.g. for MattNet $86.1 \rightarrow 51.2 \rightarrow 16.1$). (6) From the Oracle performance on COCO versus non-COCO expressions, we see that all models are able to significantly improve their performance even for non-COCO objects (e.g. for DBNet 37.3 to 59.0). DBNet benefits more than MattNet from Oracle boxes, showing its higher potential to generalize to a new domain given better proposals.

## 6    Video object segmentation results

In this section we present single and multiple video object segmentation results using natural language referring expressions on two datasets: DAVIS$_{16}$ [38] and DAVIS$_{17}$ [42]. In addition, we experiment with fusing two complementary sources of information, employing both the pixel-level mask and language supervision on the first frame. All results here are obtained using the bounding boxes given by the MattNet model [56] trained on RefCOCO [57] after the temporal consistency step (see §3.1).

For evaluation we use the IoU measure (also called Jaccard index - $J$) between the ground truth and the predicted segmentation, averaged across all video sequences and all frames. For DAVIS$_{17}$ we also employ the $J\&F$ measure proposed in [42].

### 6.1    DAVIS$_{16}$ single object segmentation

Table 3 compares our results to previous work on DAVIS$_{16}$ [38]. As we employ MattNet [56], which exploits Mask R-CNN [12] box proposals, we also would like to compare to its segments. We report the oracle Mask R-CNN results, where on each frame the segment with the highest ground truth overlap was chosen. Even with the oracle assignment of segments, [12] under-performs compared to our segmentation model (71.5 vs. 83.1). This shows that for very detailed mask annotations (as in DAVIS$_{16/17}$) a more complex segmentation module than the Mask R-CNN segmentation head is required (which itself is a shallow FCN with reduced output resolution, resulting in coarse masks).

Our method, while only exploiting language, shows competitive performance, on par with techniques which use a pixel-level mask on the first frame (82.8 vs. 81.7 for

| Supervision | AC | LR | SV | SC | CS | DB | BC | FM | MB | DEF | OCC |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Language | 80.1 | **79.0** | 74.4 | 77.6 | 85.7 | 66.4 | **85.0** | 77.7 | 78.1 | 84.3 | 80.1 |
| Mask | **81.2** | 78.1 | 75.9 | 79.0 | 85.6 | 68.0 | 82.8 | 79.0 | 79.9 | 85.6 | 80.5 |
| Mask + Lang. | 81.0 | **79.0** | **76.8** | **80.4** | **86.8** | **72.2** | 84.4 | **79.5** | **80.4** | **85.9** | **82.3** |

Table 4: Attribute-based results with different forms of supervision on DAVIS$_{16}$ val set. AC: appearance change, LR: low resolution, SV: scale variation, SC: shape complexity, CS: camera shake, DB: dynamic background, BC: background clutter, FM: fast motion, MB: motion blur, DEF: deformation, OCC: occlusions. See §6.1 for more details.

OnAVOS [50]). This shows that high quality results can be obtained via a more natural way of human-computer interaction – referring to an object via language, making video segmentation techniques more applicable in practice. Compared to mask supervision employing language results in a runtime speed up: it is $\sim 15$ times faster to specify the object with language (79s [26] vs. 5s) plus online tuning is not needed for good performance ([33] reports 10min for online tuning with 80.2 vs. our 82.8). Note that [33,50] show superior results to our approach ($\sim$ 86 mIoU). However, they employ additional cues by incorporating semantic information [33] or doing online adaptation [50]. Potentially, these techniques can also be applied to our method, though it is out of scope of this paper.

Compared to the approaches which use point click supervision [2,32], our method shows superior performance (82.8 vs. 80.6 and 80.9). This indicates that language can be successfully utilized as an alternative and cheaper form of supervision for video object segmentation, on par with clicks and scribbles.

**Maks and language.** In Table 3 we also report the results for variants using only mask supervision on the the first frame or combining both mask and language (see §3.2 for details). Notice that employing either mask or language results in comparable performance (82.8 vs. 83.1), while fusing both modalities leads to a further improvement (82.8 vs. 84.5).

| Supervision | Method | mIoU |
|---|---|---|
| Oracle | Mask R-CNN [12] | 71.5 |
| Unsupervised | FusionSeg [17] | 70.7 |
| | LVO [47] | 75.9 |
| | ARP [21] | 76.2 |
| 1st frame mask | SegFlow [6] | 76.1 |
| | MaskTrack [39] | 79.7 |
| | OSVOS[1] [33] | 80.2 |
| | MaskRNN [15] | 80.4 |
| | OnAVOS[2] [50] | 81.7 |
| | Our | 83.1 |
| Clicks | iVOS[2] | 80.6 |
| | DEXTR [32] | 80.9 |
| Language | Our | 82.8 |
| Mask + Lang. | Our | **84.5** |

(Semi-supervised)

Table 3: Comparison of video object segmentation results on DAVIS$_{16}$, val set.

This shows that referring expressions are complementary to visual forms of supervision and can be exploited as an additional source of guidance for segmentation, on top of not only pixel-level masks, but potentially scribbles and point clicks.

Table 4 presents a more detailed evaluation using video attributes. We report the averaged results on a subset of sequences where a certain challenging attribute is present. Note that using language alone leads to more robust performance for videos with low resolution, camera shake and background clutter without the need for an expensive

---

[1] OSVOS$^S$ reports 86.0 mIoU by employing semantic segmentation as additional supervision.

[2] OnAVOS gives 86.1 mIoU by exploiting online adaptation on successive frames.

pixel-level mask. When utilizing both mask and language we observe that the system becomes consistently more robust to various video challenges (e.g. fast motion, occlusions, motion blur, etc.) and compares favorably to mask only on all attributes, except appearance change. Overall, employing language can help the model to better handle occlusions, avoid drift and better adapt to complex dynamics inherent to video.

**Ablation study.** We validate the contributions of the components in our method (see §3) by presenting an ablation study in Table 5 on $DAVIS_{16}$, training set. Augmenting the ground truth boxes by random jittering makes the system more robust to sloppy boxes at test time (82.5 vs. 80.6), while employing motion cues allows to better handle moving objects

| Variant | mIoU | $\Delta$ |
|---|---|---|
| Full system | 82.5 | - |
| No box jittering | 80.6 | $-1.9$ |
| No optical flow magnitude | 75.9 | $-4.7$ |
| No temporal consistency | 72.5 | $-3.4$ |
| Backbone architecture of [39] | 72.2 | $-3.7$ |

Table 5: Ablation study on $DAVIS_{16}$.

(80.6 vs. 75.9). Temporal consistency step helps to provide more temporally coherent boxes (4.3 mIoU point boost for grounding, see Table 1) and hence improve the final segmentation quality (75.9 vs. 72.5). Exploiting the proposed network architecture versus using the network proposed in [39] results in 3.7 point boost (75.9 vs. 72.2), providing more detailed object masks. Overall, all components introduced in our approach lead to the state-of-the-art results on $DAVIS_{16}$.

### 6.2   $DAVIS_{17}$ multiple object segmentation

Table 6 presents results on $DAVIS_{17}$ [42]. The lower numbers in comparison with Table 3 indicate that $DAVIS_{17}$ is significantly more difficult than $DAVIS_{16}$. Even when employing mask supervision on the first frame the dataset presents a challenging task and there is much room for improvement. The semi-supervised methods perform well on foreground-background segmentation, but have problems separating multiple foreground objects, handling small objects and preserving the correct object identities [42].

Compared to mask supervision using language descriptions significantly underperforms. We believe that one of the main problems is a relatively unstable behavior of the underlying grounding model. There are a lot of identity switches, that are heavily penalized by the evaluation metric as every pixel should be assigned to one instance. We conducted an oracle experiment assigning Mask R-CNN box proposals to the correct object ids and then performing segmentation (denoted "Oracle - Grounding"). We observe a significant increase in performance (37.3 to 54.9), making the results competitive to mask supervision. If we utilize Mask R-CNN segment proposals for oracle case, the result is 2.1 points lower than using our segmentation model on top. The underlying choice of proposals for the grounding model could also have its effect. If the object is not detected by Mask R-CNN, the grounding model has no chances to recover the correct instance. To evaluate the influence of proposals we conduct an oracle experiment where the ground truth boxes are exploited in the grounding model (denoted "Oracle - Box proposals"). With oracle boxes we observe an increase in performance (37.3 to 42.1), however, recovering the correct identities still poses a problem for grounding.

Another factor influencing the results is the domain shift between the training and test data. Both Mask R-CNN and MattNet are trained on MS COCO [26], and have troubles recovering instances not belonging to 80 COCO categories. We split the $DAVIS_{17}$ validation set into COCO and non-COCO objects/language queries (43 vs. 18) and eval-

uate separately on two subsets. As in §5, we observe much higher results for COCO queries (45 to 27.5), indicating the problem of generalization from training to test data.

The method which exploits scribble supervision [41] performs on par with our approach. Note that even for scribble supervision the task remains difficult.

**Mask and language.** In Table 6 we also report the results for variants of our approach using only mask supervision or combining mask and language. Employing language on top of mask leads to an increase in performance over using mask only (58 to 59), again showing complementarity of both sources of supervision.

Figure A1 provides qualitative results of our method using only language as supervision. We observe successful handling of similar looking objects, fast motion, deformations and partial occlusions.

**Discussion.** Our results indicate that language alone can be successfully used as an alternative and a more natural form of supervision. Particularly, high quality results can be achieved for videos with the salient target object. Videos with

| Supervision | Method | mIoU | $J\&F$ |
|---|---|---|---|
| Oracle | Mask R-CNN [12] | 52.8 | 53.3 |
|  | Grounding | 54.9 | 57.4 |
|  | Box proposals | 42.1 | 45.3 |
| 1st frame mask | OSVOS [3] | 52.1 | 57.0 |
|  | OnAVOS[3] [49] | 57.0 | 59.4 |
|  | MaskRNN [15] | 60.5 | - |
|  | Our | 58.0 | 60.8 |
| Scribbles | CNN lin. class. [41] | - | 39.3 |
|  | Scribble-OSVOS [41] | - | 39.9 |
| Language | Our | 37.3 | 39.3 |
|  | Our, COCO | *45.0* | *47.5* |
|  | Our, non-C. | *27.5* | *29.4* |
| Mask+Lang. | Our | 59.0 | 62.2 |

Table 6: Comparison of semi-supervised video object segmentation methods on DAVIS$_{17}$, val set. Numbers in italic are reported on subsets of DAVIS$_{17}$ containing/non-containing COCO objects.

multiple similar looking objects pose a challenge for grounding models, as they have problems preserving object identities across frames. Experimentally we show that better proposals, grounding and proximity of training and test data can further boost the performance for videos with multiple objects. Language is complementary to mask supervision and can be exploited as an additional source of guidance for segmentation.

## 7  Conclusion

In this work we propose the task of video object segmentation using language referring expressions. We propose an approach to address this new task as well as extend two well-known video object segmentation benchmarks with textual descriptions of target objects. Our experiments indicate that language alone can be successfully exploited to obtain high quality segmentations of objects in videos. While allowing a more natural human-computer interaction, using guidance from language descriptions can also make video segmentation more robust to occlusions, complex dynamics and cluttered backgrounds. We show that classical semi-supervised video object segmentation which uses the mask annotation on the first frame can be further improved by the use of language descriptions. We believe there is a lot of potential in fusing lingual (referring expressions) and visual (clicks, scribbles or masks) forms of supervision for object segmentation in video. We hope that our results encourage more research on video object segmentation with referring expressions and foster discovery of new techniques applicable in realistic settings, which discard tedious pixel-level annotations.

---

[3] OnAVOS reports 64.5 mIoU by performing online adaptation on successive frames.

# References

1. Balajee Vasudevan, A., Dai, D., Van Gool, L.: Object referring in videos with language and human gaze. In: CVPR (2018)
2. Benard, A., Gygli, M.: Interactive video object segmentation in the wild. arXiv: 1801.00269 (2017)
3. Caelles, S., Maninis, K.K., Pont-Tuset, J., Leal-Taixe, L., Cremers, D., Gool, L.V.: One-shot video object segmentation. In: CVPR (2017)
4. Chen, K., Kovvuri, R., Nevatia, R.: Query-guided regression network with context policy for phrase grounding. In: ICCV (2017)
5. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. arXiv:1606.00915 (2016)
6. Cheng, J., Tsai, Y.H., Wang, S., Yang, M.H.: Segflow: Joint learning for video object segmentation and optical flow. In: ICCV (2017)
7. Cheng, M.M., Mitra, N.J., Huang, X., Torr, P.H.S., Hu, S.M.: Global contrast based salient region detection. PAMI (2015)
8. Dollár, P., Zitnick, C.L.: Fast edge detection using structured forests. PAMI (2015)
9. Faktor, A., Irani, M.: Video segmentation by non-local consensus voting. In: BMVC (2014)
10. Gavrilyuk, K., Ghodrati, A., Li, Z., Snoek, C.G.: Actor and action video segmentation from a sentence. In: CVPR (2018)
11. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: AISTATS (2010)
12. He, K., Gkioxari, G., Dollár, P., Girshick, R.B.: Mask R-CNN. In: ICCV (2017)
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
14. Hu, R., Rohrbach, M., Andreas, J., Darrell, T., Saenko, K.: Modeling relationships in referential expressions with compositional modular networks. In: CVPR (2017)
15. Hu, Y.T., Huang, J., Schwing, A.G.: Maskrnn: Instance level video object segmentation. In: NIPS (2017)
16. Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., Brox, T.: Flownet 2.0: Evolution of optical flow estimation with deep networks. In: CVPR (2017)
17. Jain, S.D., Xiong, B., Grauman, K.: Fusionseg: Learning to combine motion and appearance for fully automatic segmention of generic objects in videos. In: CVPR (2017)
18. Jampani, V., Gadde, R., Gehler, P.V.: Video propagation networks. arXiv:1612.05478 (2016)
19. Jang, W.D., Kim, C.S.: Online video object segmentation via convolutional trident network. In: CVPR (2017)
20. Khoreva, A., Benenson, R., Ilg, E., Brox, T., Schiele, B.: Lucid data dreaming for multiple object tracking. arXiv: 1703.09554 (2017)
21. Koh, Y., Kim, C.: Primary object segmentation in videos based on region augmentation and reduction. In: CVPR (2017)
22. Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L.J., Shamma, D.A., Bernstein, M., Fei-Fei, L.: Visual genome: Connecting language and vision using crowdsourced dense image annotations. arXiv: 1602.07332 (2016)
23. Li, R., Li, K., Kuo, Y., Shu, M., Qi, X., Shen, X., Jia, J.: Referring image segmentation via recurrent refinement networks. In: CVPR (2018)
24. Li, Z., Tao, R., Gavves, E., Snoek, C.G.M., Smeulders, A.W.M.: Tracking by natural language specification. In: CVPR (2017)
25. Lin, D., Dai, J., Jia, J., He, K., Sun, J.: Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In: CVPR (2016)

26. Lin, T., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014)
27. Liu, C., Lin, Z., Shen, X., Yang, J., Lu, X., Yuille, A.: Recurrent multimodal interaction for referring image segmentation. In: ICCV (2017)
28. Liu, W., Rabinovich, A., Berg, A.C.: Parsenet: Looking wider to see better. arxiv:1506.04579 (2015)
29. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: CVPR (2015)
30. Luo, R., Shakhnarovich, G.: Comprehension-guided referring expressions. In: CVPR (2017)
31. Maerki, N., Perazzi, F., Wang, O., Sorkine-Hornung, A.: Bilateral space video segmentation. In: CVPR (2016)
32. Maninis, K., Caelles, S., Pont-Tuset, J., Gool, L.V.: Deep extreme cut: From extreme points to object segmentation. In: CVPR (2018)
33. Maninis, K., Caelles, S., Chen, Y., Pont-Tuset, J., Leal-Taixé, L., Cremers, D., Gool, L.V.: Video object segmentation without temporal information. arxiv: 1709.06031 (2017)
34. Mao, J., Jonathan, H., Toshev, A., Camburu, O., Yuille, A., Murphy, K.: Generation and comprehension of unambiguous object descriptions. In: CVPR (2016)
35. Nagaraja, V.K., Morariu, V.I., Davis, L.S.: Modeling context between objects for referring expression understanding. In: ECCV (2016)
36. Papazoglou, A., Ferrari, V.: Fast object segmentation in unconstrained video. In: ICCV (2013)
37. Papazoglou, A., Ferrari, V.: Fast object segmentation in unconstrained video. In: ICCV (2013)
38. Perazzi, F., Pont-Tuset, J., McWilliams, B., Gool, L.V., Gross, M., Sorkine-Hornung, A.: A benchmark dataset and evaluation methodology for video object segmentation. In: CVPR (2016)
39. Perazzi, F., Khoreva, A., Benenson, R., Schiele, B., Sorkine-Hornung, A.: Learning video object segmentation from static images. In: CVPR (2017)
40. Perazzi, F., Wang, O., Gross, M., Sorkine-Hornung, A.: Fully connected object proposals for video segmentation. In: ICCV (2015)
41. Pont-Tuset, J., Caelles, S., Perazzi, F., Montes, A., Maninis, K.K., Chen, Y., Van Gool, L.: The 2018 davis challenge on video object segmentation. arXiv:1803.00557 (2018)
42. Pont-Tuset, J., Perazzi, F., Caelles, S., Arbeláez, P., Sorkine-Hornung, A., Van Gool, L.: The 2017 davis challenge on video object segmentation. arXiv:1704.00675 (2017)
43. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. IJCV (2015)
44. Shin Yoon, J., Rameau, F., Kim, J., Lee, S., Shin, S., So Kweon, I.: Pixel-level matching for video object segmentation using convolutional neural networks. In: ICCV (2017)
45. Thomee, B., Shamma, D.A., Friedland, G., Elizalde, B., Ni, K., Poland, D., Borth, D., Li, L.J.: Yfcc100m: the new data in multimedia research. Communications of the ACM (2016)
46. Tokmakov, P., Alahari, K., Schmid, C.: Learning motion patterns in videos. In: CVPR (2017)
47. Tokmakov, P., Alahari, K., Schmid, C.: Learning video object segmentation with visual memory. In: ICCV (2017)
48. Tsai, Y.H., Yang, M.H., Black, M.J.: Video segmentation via object flow. In: CVPR (2016)
49. Voigtlaender, P., Leibe, B.: Online adaptation of convolutional neural networks for the 2017 davis challenge on video object segmentation. DAVIS Challenge - CVPR Workshops (2017)
50. Voigtlaender, P., Leibe, B.: Online adaptation of convolutional neural networks for video object segmentation. In: BMVC (2017)
51. Wang, P., Chen, P., Yuan, Y., Liu, D., Huang, Z., Hou, X., Cottrell, G.: Understanding convolution for semantic segmentation. arXiv:1702.08502 (2017)

52. Wang, W., Shen, J.: Super-trajectory for video segmentation. arXiv:1702.08634 (2017)
53. Wen, L., Du, D., Lei, Z., Li, S.Z., Yang, M.H.: Jots: Joint online tracking and segmentation. In: CVPR (2015)
54. Xiao, F., Lee, Y.J.: Track and segment: An iterative unsupervised approach for video object proposals. In: CVPR (2016)
55. Yeh, R., Xiong, J., Hwu, W.M., Do, M., Schwing, A.: Interpretable and globally optimal prediction for textual grounding using image concepts. In: NIPS (2017)
56. Yu, L., Lin, Z., Shen, X., Yang, J., Lu, X., Bansal, M., Berg, T.L.: Mattnet: Modular attention network for referring expression comprehension. In: CVPR (2018)
57. Yu, L., Poirson, P., Yang, S., Berg, A.C., Berg, T.L.: Modeling context in referring expressions. In: ECCV (2016)
58. Zhang, Y., Yuan, L., Guo, Y., He, Z., Huang, I.A., Lee, H.: Discriminative bimodal networks for visual localization and detection with natural language queries. In: CVPR (2017)

## Supplementary Material

This supplementary material provides additional quantitative and qualitative results and is structured as follows.

Section A discusses two types of referring expressions - 1st frame vs. full video - and the effect of 1st frame annotations being invalid for the whole video. It also provides additional examples of the collected referring expressions for video object segmentation task (see Figure A2).

Section B provides additional evaluation of natural language grounding models on the Lingual ImageNet Videos [24] and compares results with the work of [24] (Table B1).

Section C provides additional evaluation metrics for DAVIS$_{16}$ (Table C2) and comparisons of different grounding models, effect of temporal consistency and annotation types on video object segmentation task (Table C3). We also include more qualitative examples for Language, Mask and Mask + Language approaches (see Figures C3-C5).

## A      Referring expressions for video object segmentation

As our goal is to segment objects in videos using language specifications, we augment all objects annotated with mask labels in DAVIS$_{16}$ [38] and DAVIS$_{17}$ [42] with non-amb-iguous referring expressions.

We collected referring expression annotations using two different settings, asking the annotators to provide a description of the target object based on the first frame only as well as on the full video. Future work may choose which setting they prefer more.

We experiment with both annotation types. While the first type is more similar to image-based referring expressions, the



*Original query: "A brown camel" vs.*

*Corrected: "A brown camel in the front"*

Figure A1:  Predictions for the ambiguous query and its correction.

second type has different trends, tending to be more complex/long due to increased complexity of the video. We report the grounding (Table 1 in the main paper) and VOS results (Table C3) with both types, showing that DBNet [58] benefits from the "full video" descriptions, while MattNet [57] has difficulties coping with more complex language.

Concerned that the referring expressions obtained by only looking at the 1st frame might be potentially invalid for the entire video, on DAVIS$_{17}$ we ask a user to mark which 1st frame expressions become ambiguous/invalid over time, and to correct them to be valid for the full video (e.g. Fig A1). Only ∼15% of all descriptions were marked invalid. Though some descriptions become ambiguous/invalid over time, it does not impact strongly the results (original 36.9 vs. corrected 37.1 mIoU). One of the reasons is that *temporal consistency* helps to disambiguate some of such cases (Fig A1). Another reason is that invalid descriptions might still contain valid info (e.g. "a boy in red on the left", the boy is no longer on the left, but still in red).

We present additional examples of collected referring expressions in Figure A2.
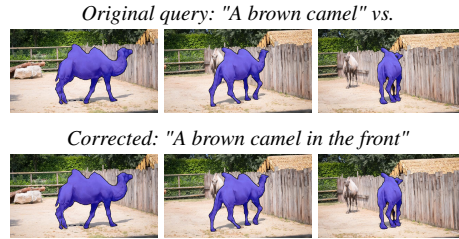
ID 1: "A man on the left wearing blue"
ID 2: "A man on the right wearing red"
ID 3: "A referee in the middle in white"

ID 1: "A man in a blue dress on the left getting punched"
ID 2: "A man in a red dress on the right punching"
ID 3: "A man in a white shirt and black shorts in the middle"

ID 1: "A brown sheep in the middle"
ID 2: "A sheep on the left with a black face"
ID 3: "A black lamb with white nose"
ID 4: "A white lamb next to a brown sheep"
ID 5: "A white lamb in the middle next to a white sheep"

ID 1: "A brown sheep in the front"
ID 2: "A grey sheep with dark face moving behind fence"
ID 3: "A black baby sheep"
ID 4: "A white baby sheep closer to a brown sheep"
ID 5: "A white baby sheep farther from a brown sheep"

ID 1: "A black bicycle"
ID 2: "A backpack"
ID 3: "A black board"
ID 4: "A man on a bicycle in a black jacket"
ID 5: "A man in a yellow t-shirt"

ID 1: "A bicycle moving on the road"
ID 2: "A backpack worn by a guy"
ID 3: "A longboard"
ID 4: "A guy riding a bicycle"
ID 5: "A person rolling over longboard"

First frame annotation                    Full video annotation

Figure A2: Example of collected annotations provided for the first frame (left) vs. the full video (right). Full video annotations include descriptions of activities and overall are more complex than the ones provided for the first frame.

## B  Language grounding results on Lingual ImageNet Videos

For the natural language grounding task we additionally consider Lingual ImageNet Videos [24], which provides referring expression annotations for a subset of the ImageNet Video Object Detection dataset [43]. The dataset is split into a training and a validation set, each consisting of 50 videos. The performance on Lingual ImageNet [24] is measured in terms of the AUC (area under the curve) score metric, following [24].

| Method | Supervision | AUC score |
|---|---|---|
| Tracking by language [24] | Language | 26.3 |
| | Box | 47.9 |
| | Box + Language | 49.4 |
| DBNet | Language | 54.0 |
| MattNet | Language | **60.8** |

Table B1: Comparison of grounding models on Lingual ImageNet Videos, val set.

Here we compare to [24], who perform tracking of objects using language specifications. Table B1 presents grounding results reported by [24], including tracking by language only, tracking given the ground-truth bounding box on the first frame, and the combined approach. Our method is based on language input only, specifically, we report the results after the temporal consistency step applied to DBNet and MattNet

| Supervision | Method | DAVIS$_{16}$ | | | | | | Temp. stab., $T$ |
| | | Region, $J$ | | | Boundary, $F$ | | | |
| | | Mean ↑ | Recall ↑ | Decay ↓ | Mean ↑ | Recall ↑ | Decay ↓ | Mean ↓ |
|---|---|---|---|---|---|---|---|---|
| Oracle | Mask R-CNN [12] | 71.5 | 87.3 | 5.9 | 72.4 | 84.6 | 6.8 | 24.8 |
| Unsupervised | NLC [9] | 55.1 | 55.8 | 12.6 | 52.3 | 51.9 | 11.4 | 42.5 |
| | FST[36] | 55.8 | 64.9 | 0.0 | 51.1 | 51.6 | 2.9 | 36.6 |
| | SegFlow[6] | 67.4 | 81.4 | 6.2 | 66.7 | 77.1 | 5.1 | 28.2 |
| | MP-Net [46] | 70.0 | 85.0 | 1.3 | 65.9 | 79.2 | 2.5 | 57.2 |
| | FusionSeg [17] | 70.7 | 83.5 | 1.5 | 65.3 | 73.8 | 1.8 | 32.8 |
| | LVO [47] | 75.9 | 89.1 | 0.0 | 72.1 | 8.4 | 1.3 | 26.5 |
| | ARP [21] | 76.2 | 91.1 | 7.0 | 70.6 | 83.5 | 7.9 | 39.3 |
| Semi-supervised / 1st frame mask | FCP [40] | 58.4 | 71.5 | **-2.0** | 49.2 | 49.5 | **-1.1** | 30.6 |
| | BVS [31] | 60.0 | 66.9 | 28.9 | 58.8 | 67.9 | 21.3 | 34.7 |
| | ObjFlow [48] | 68.0 | 75.6 | 26.4 | 63.4 | 70.4 | 27.2 | 22.2 |
| | PLM [44] | 70.2 | 86.3 | 11.2 | 62.5 | 73.2 | 14.7 | 31.8 |
| | VPN [18] | 70.2 | 82.3 | 12.4 | 65.5 | 69.0 | 14.4 | 32.4 |
| | CTN [19] | 73.5 | 87.4 | 15.6 | 69.3 | 79.6 | 12.9 | 22.0 |
| | SegFlow [6] | 76.1 | 90.6 | 12.1 | 76.0 | 85.5 | 10.4 | **18.9** |
| | MaskTrack [39] | 79.7 | 93.1 | 8.9 | 75.4 | 87.1 | 9.0 | 21.8 |
| | OSVOS [3] | 79.8 | 93.6 | 14.9 | 80.6 | 92.6 | 15.0 | 37.8 |
| | MaskRNN [15] | 80.4 | 96.0 | 4.4 | 82.3 | 93.2 | 8.8 | 19.0 |
| | OnAVOS[1] [50] | 81.7 | 92.2 | 11.9 | 81.1 | 88.2 | 11.2 | 27.3 |
| | Our | 83.1 | 95.1 | 9.8 | 85.7 | 94.4 | 9.6 | 24.0 |
| Language | Our | 82.8 | 94.1 | 3.2 | 85.4 | 94.7 | 3.4 | 22.6 |
| Mask + Lang. | Our | **84.5** | **96.3** | 8.2 | **86.9** | **95.9** | 8.7 | 24.8 |

Table C2: Comparison of video object segmentation results on DAVIS$_{16}$, validation set.

predictions. As we see both models significantly outperform [24], even when [24] has access to the ground-truth bounding box on the first frame.

## C    Video object segmentation

### C.1    Additional metrics for DAVIS$_{16}$

We report video object segmentation results for the DAVIS$_{16}$ benchmark in Table C2, using evaluation metrics proposed in [38]. Three measures are used: region similarity in terms of intersection-over-union ($J$, higher is better), contour accuracy ($F$, higher is better), and temporal instability of the masks ($T$, lower is better). See [38] for more details. Note that using only language supervision results in a smaller decay over time for $J$ and $F$ measures and a better overall temporal stability $T$ compared to employing pixel-level mask supervision on the first frame.

---

[1] OnAVOS gives 86.1 mIoU by online adaptation on successive frames.

| Annotation type | Grounding | Temporal consistency | mIoU | $J\&F$ |
|---|---|---|---|---|
| 1st frame | DBNet | - | 32.6 | 34.7 |
| | | ✓ | 35.4 | 37.6 |
| | MattNet | - | 35.4 | 38.5 |
| | | ✓ | **37.3** | **39.3** |
| Full video | DBNet | ✓ | 35.5 | 37.7 |
| | MattNet | ✓ | 35.5 | 37.1 |

Table C3: Effect of different grounding models, temporal consistency and annotation types on video object segmentation on DAVIS$_{17}$, validation set.

## C.2    Effect of grounding models, temporal consistency and annotation types on video object segmentation

Table C3 reports the effect of different grounding models, temporal consistency step for grounding and employing the first frame versus the full video descriptions on video object segmentation.

We compare DBNet versus MattNet (trained on RefCOCO [57]) as a base grounding model for video object segmentation task. Exploiting MattNet grounding boxes results in a better performance compared to DBNet (37.3 vs. 35.4). Overall the temporal consistency step has a positive impact on video object segmentation performance across different grounding models (for MattNet $35.4 \rightarrow 37.3$ and for DBNet $32.6 \rightarrow 35.4$).

We also compare the segmentation performance from first frame versus full video descriptions in Table C3. Employing the full video versus the first frame descriptions results in a minor improvement for DBNet (35.4 vs. 35.5), however has a negative effect for MattNet (37.3 vs. 35.5). The same diverging has been observed for language grounding results in the main paper when comparing results on expressions provided for the first frame versus expressions provided for the full video in Table 2. We attribute this to the fact that DBNet is trained on the more diverse Visual Genome descriptions and can handle better more complex full video expressions.

## C.3    Qualitative results for video object segmentation

Figure C3 provides more qualitative examples of Language-only supervision for video object segmentation on DAVIS$_{16}$ and DAVIS$_{17}$, validation sets. We observe successful handling of shape deformations, fast motion as well as partial and full occlusions.

Figure C4 shows examples of Mask + Language supervision on DAVIS$_{17}$, validation set. We observe high quality instance level segmentation of multiple similar looking objects.

Figure C5 shows comparison of Language versus Mask supervision on DAVIS$_{16}$ and DAVIS$_{17}$, validation sets. Note that using only language supervision results in a more robust performance for videos with similar looking instances and camera view changes in comparison to employing pixel-level masks.

*ID 1: "A red car".*

*ID 1: "A man jumping across fences".*

*ID 1: "A dog running in the garden".*

*ID 1: "A goat walking on rocks".*

*ID 1: "A red and white car".*

*ID 1: "A woman riding a horse". ID 2: "A horse doing high-jumps".*

*ID 1: "A bald man with black belt in the center". ID 2: "A man with blue belt on the right".*

*ID 1: "A boy wearing a white t-shirt". ID 2: "A red bmx bike".*

*ID 1: "A green motorbike". ID 2: "A man riding a motorbike".*

Figure C3: Video object segmentation qualitative results using only Language as supervision on DAVIS$_{16}$ and DAVIS$_{17}$, val sets. Frames sampled along the video duration.

*ID 1: "A man wearing a cap". ID 2: "A black bike".*



*ID 1: "A brown piglet in the middle". ID 2: "A brown and white colored piglet".*
*ID 3: "An adult pig on the right".*



*ID 1: "An orange goldfish in the center next to the largest fish". ID 2: "The biggest goldfish".*
*ID 3: "The smallest goldfish". ID 4: "A small goldfish in the end".*
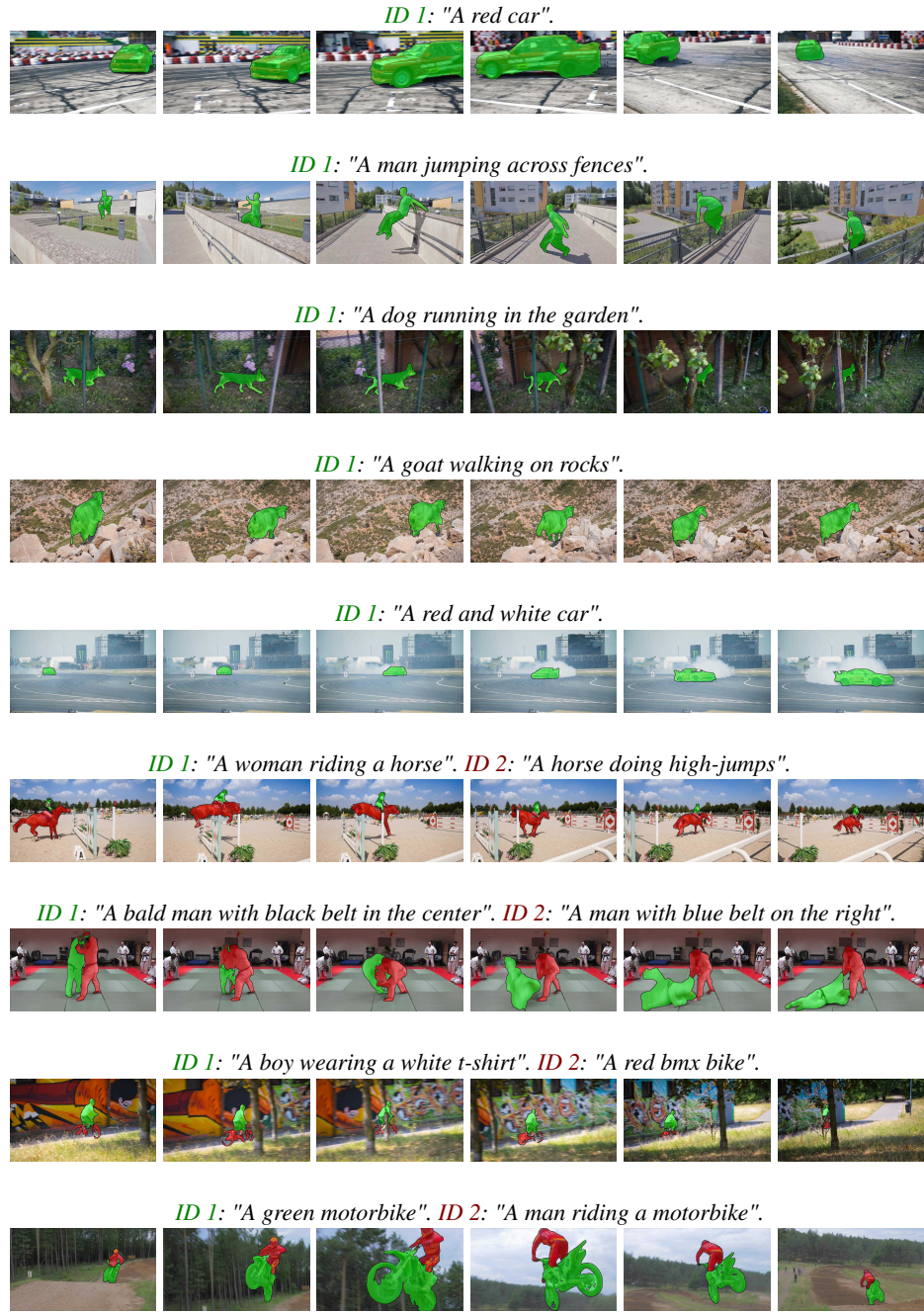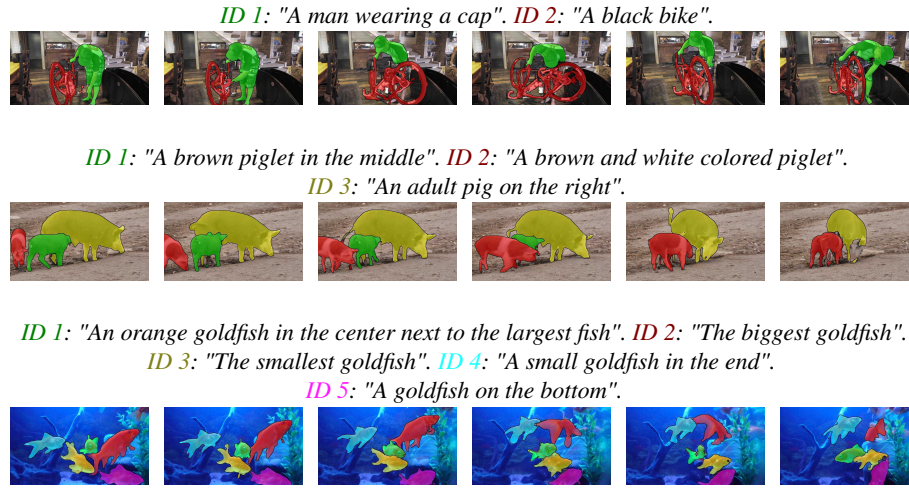*ID 5: "A goldfish on the bottom".*



Figure C4: Video object segmentation qualitative results using Mask + Language as supervision on DAVIS$_{17}$, val set. Frames sampled along the video duration. In the last row we visualize a failure case of the proposed approach.

Language supervision, *ID 1: "A brown camel in the front".*



Pixel-level mask supervision



Language supervision, *ID 1: "A silver car".*



Pixel-level mask supervision



Language supervision, *ID 1: "A black car".*



Pixel-level mask supervision



Language supervision, *ID 1: "A green motorbike". ID 2: "A man riding a motorbike".*



Pixel-level mask supervision



Language supervision, *ID 1: "A black scooter ridden by a man".*
*ID 2: "A man in a suit riding a scooter".*



Pixel-level mask supervision



Figure C5: Video object segmentation results using Language versus Mask on the 1st frame as supervision on DAVIS$_{16}$ and DAVIS$_{17}$, val sets. Using language only results in a more robust performance for videos with similar looking instances and camera view changes in comparison to employing pixel-level masks. Frames sampled along the video duration. The videos are chosen with the highest mIoU difference.