

Video Object Segmentation Without Temporal Information

K.-K. Maninis*, S. Caelles*, Y. Chen,
J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool

Abstract—Video Object Segmentation, and video processing in general, has been historically dominated by methods that rely on the temporal consistency and redundancy in consecutive video frames. When the temporal smoothness is suddenly broken, such as when an object is occluded, or some frames are missing in a sequence, the result of these methods can deteriorate significantly. This paper explores the orthogonal approach of processing each frame independently, i.e. disregarding the temporal information. In particular, it tackles the task of semi-supervised video object segmentation: the separation of an object from the background in a video, given its mask in the first frame. We present Semantic One-Shot Video Object Segmentation (OSVOS^S), based on a fully-convolutional neural network architecture that is able to successively transfer generic semantic information, learned on ImageNet, to the task of foreground segmentation, and finally to learning the appearance of a single annotated object of the test sequence (hence one shot). We show that instance-level semantic information, when combined effectively, can dramatically improve the results of our previous method, OSVOS. We perform experiments on two recent single-object video segmentation databases, which show that OSVOS^S is both the fastest and most accurate method in the state of the art. Experiments on multi-object video segmentation show that OSVOS^S obtains competitive results.

Index Terms—Video Object Segmentation, Convolutional Neural Networks, Semantic Segmentation, Instance Segmentation.



1 INTRODUCTION

A video is a temporal sequence of static images that give the impression of continuous motion when played consecutively and rapidly. The illusion of motion pictures is due to the persistence of human vision [23], [65], [71]: the fact that it cannot perceive very high frequency changes [71] because of the temporal integration of incoming light into the retina [65]. This property has been exploited since the appearance of the phenakistoscope [63] or the zoetrope [23], which displayed a sequence of drawings creating the illusion of continuous movement.

In order to achieve the high frequency to produce the video illusion, consecutive images vary very smoothly and slowly: the information in a video is very redundant and neighboring frames carry very similar information. In video coding, for instance, this is the key idea behind video compression algorithms such as motion-compensated coding [65], where instead of storing each frame independently, one picks a certain image and only codes the modifications to be done to it to generate the next frame.

Video processing in general, and video segmentation in particular, is also dominated by this idea, where *motion estimation* has emerged as a key ingredient for some of the state-of-the-art video segmentation algorithms [17], [27], [49], [56], [67]. Exploiting it is not a trivial task however, as one has to compute temporal matches in the form of optical flow or dense trajectories [4], which can be an even harder problem to solve.

On the other hand, processing each frame independently would allow us to easily parallelize the computation, and to not be affected by sequence interruptions, to process the frames at any desire rate, etc. This paper explores how to segment objects in

videos when processing each frame independently, that is, by ignoring the temporal information and redundancy. In other words, we cast video object segmentation as a per-frame segmentation problem given the *model* of the object from one (or various) manually-segmented frames.

This stands in contrast to the dominant approach where temporal consistency plays the central role, assuming that objects do not change too much between one frame and the next. Such methods adapt their single-frame models smoothly throughout the video, looking for targets whose shape and appearance vary *gradually* in consecutive frames, but fail when those constraints do not apply, unable to recover from relatively common situations such as occlusions and abrupt motion.

We argue that temporal consistency was needed in the past, as one had to overcome major drawbacks of the then inaccurate shape or appearance models. On the other hand, in this paper deep learning will be shown to provide a sufficiently accurate model of the target object to produce very accurate results even when processing each frame independently. This has some natural advantages: OSVOS^S is able to segment objects throughout occlusions, it is not limited to certain ranges of motion, it does not need to process frames sequentially, and errors are not temporally propagated. In practice, this allows OSVOS^S to handle e.g. interlaced videos of surveillance scenarios, where cameras can go blind for a while before coming back on again.

Given the first frame, we create an *appearance model* of the object of interest and then look for the pixels that better match this model in the rest of the frames. To do so, we will make use of Convolutional Neural Networks (CNNs), which are revolutionizing many fields of computer vision. For instance, they have dramatically boosted the performance for problems like image classification [22], [32], [62] and object detection [15], [16], [39]. Image segmentation has also been taken over by CNNs

- K.-K. Maninis, S. Caelles, Y. Chen, J. Pont-Tuset, and L. Van Gool are with the ETHZ, Zürich. First two authors contributed equally.
- L. Leal-Taixé and D. Cremers are with the TUM, München.
- Contacts in <http://www.vision.ee.ethz.ch/~cvlsegmentation/>

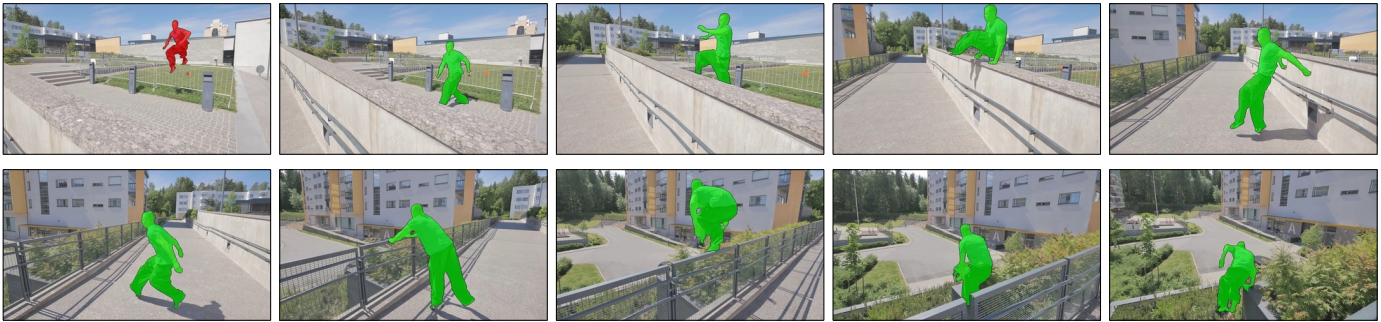


Fig. 1. **Example result of our technique:** The segmentation of the first frame (red) is used to learn the model of the specific object to track, which is segmented in the rest of the frames independently (green). One every 10 frames shown of 90 in total.

recently [2], [3], [31], [41], [72], with deep architectures pre-trained on the weakly related task of image classification on ImageNet [60]. One of the major downsides of deep network approaches, however, is their hunger for training data. Yet, with various pre-trained network architectures one may ask how much training data do we really need for the specific problem at hand? This paper investigates segmenting an object along an entire video, when we only have one single labeled training example, e.g. the first frame.

Figure 1 shows an example result of OSVOS^S, where the input is the segmentation of the first frame (in red), and the output is the mask of the object in the 90 frames of the sequence (in green).

The first contribution of the paper is to adapt the CNN to a particular object instance given a single annotated image. To do so, we gradually adapt a CNN pre-trained on image recognition [60] to video object segmentation. This is achieved by training it on a set of videos with manually segmented objects. Finally, it is fine-tuned *at test time* on a specific object that is manually segmented in a single frame. Figure 2 shows the overview of the method. Our proposal tallies with the observation that leveraging these different levels of information to perform object segmentation would stand to reason: from generic information of a large amount of categories, passing through the knowledge of the *usual* shapes of objects in videos, down to the specific properties of a particular object we are interested in segmenting.

Our second contribution is to extend the model of the object with explicit semantic information. In the example of Figure 1, for instance, we would like to leverage the fact that we are segmenting an object of the category *person* and that there is a *single instance* of it.

In particular, we will use an *instance-aware semantic segmentation algorithm* [9], [21], [36] to extract a list of proposal of object masks in each frame, along with their categories. Given the first annotated frame, we will *infer* the categories of the objects of interest by finding the best-overlapping masks. We refer to this step as “semantic selection.”

Our method uses the extracted semantic information from the first frame to segment the rest of the video. It enforces the resulting masks to align well with the same categories selected in the first frame. If we were segmenting a person on a motorbike, then this information should be kept throughout the video. In particular, we find instances extracted from the semantic instance segmentation algorithm that best match the model of the object, and we effectively combine them with the appearance model of the object, using a conditional classifier. We call this step “semantic propagation.”

Our third contribution is that OSVOS^S can work at various points of the trade-off between speed and accuracy. In this sense, given one annotated frame, the user can choose the level of fine-tuning performed on it, giving them the freedom between a faster method or more accurate results. Experimentally, we show that OSVOS^S can run at 300 milliseconds per frame and 75.1% accuracy, and up to 86.5% when processing each frame in 4.5 seconds, for an image of 480×854 pixels.

Technically, we adopt the architecture of Fully Convolutional Networks (FCN) [14], [40], suitable for dense predictions. FCNs have recently become popular due to their performance both in terms of accuracy and computational efficiency [7], [10], [40]. Arguably, the Achilles’ heel of FCNs when it comes to segmentation is the coarse scale of the deeper layers, which leads to inaccurately localized predictions. To overcome this, a large variety of works from different fields use skip connections of larger feature maps [20], [40], [43], [72], or learnable filters to improve upscaling [47], [74].

We perform experiments on two video object segmentation datasets (DAVIS 2016 [50] and Youtube-Objects [25], [55]) and show that OSVOS^S significantly improves the state of the art in them, both in terms of accuracy and speed. We perform additional experiments for multi-object video segmentation on DAVIS 2017 [54], where we obtain competitive results by directly applying our method without adaptation to the new problem.

All resources of this paper, including training and testing code, pre-computed results, and pre-trained models will be made publicly available.

2 RELATED WORK

Semi-supervised Video Object Segmentation: Most of the current literature on semi-supervised video object segmentation enforces temporal consistency in video sequences to propagate the initial mask into the following frames. The most recent works heavily rely on optical flow, and make use of CNNs to learn to refine the mask of the object at frame n to frame $n + 1$ [28], [49] or combine the training of a CNN with ideas of bilateral filtering between consecutive frames [26]. Also, [70] follows up with the idea introduced in OSVOS and uses the result on the predicted frames on the whole sequence to further train the network at test time. Previously, and in order to reduce the computational complexity, some works make use of superpixels [6], [17], patches [13], [56], object proposals [51], or the bilateral space [46]. After that, an optimization using one of the previous aggregations of pixels is usually performed; which can

consider the full video sequence [46], [51], a subset of frames [17], or only the results in frame n to obtain the mask in $n+1$ [6], [13], [56]. As part of their pipeline, some of the methods include the computation of optical flow [17], [49], [56], or/and Conditional Random Fields (CRFs) [49] which can considerably reduce their speed. Different from those approaches, OSVOS^S is a simpler pipeline which segments each frame independently, and produces more accurate results, while also being significantly faster.

FCNs for Segmentation: Segmentation research has closely followed the innovative ideas of CNNs in the last few years. The advances observed in image recognition [22], [32], [62] have been beneficial to segmentation in many forms (semantic [40], [47], instance-level [7], [15], [52], biomedical [59], generic [41], etc.). Many of the current best performing methods are based on a deep CNN architecture, usually pre-trained on ImageNet [60], trainable end-to-end. The idea of dense predictions with CNNs was pioneered by [14] and formulated by [40] in the form of Fully Convolutional Networks (FCNs) for semantic segmentation. The authors noticed that by changing the last fully connected layers to 1×1 convolutions it is possible to train on images of arbitrary size, by predicting correspondingly-sized outputs. Their approach boosts efficiency over patch-based approaches where one needs to perform redundant computations in overlapping patches. More importantly, by removing the parameter-intensive fully connected layers, the number of trainable parameters drops significantly, facilitating training with relatively fewer labeled data.

In most CNN architectures [22], [32], [62], activations of the intermediate layers gradually decrease in size, because of spatial pooling operations or convolutions with a stride. Making dense predictions from downsampled activations results in coarsely localized outputs [40]. Deconvolutional layers that learn how to upsample are used in [47], [74] to recover accurately localized predictions. In [52], activations from shallow layers are gradually injected into the prediction to favor localization. However, these architectures come with many more trainable parameters and their use is limited to cases with sufficient data.

Following the ideas of FCNs, Xie and Tu [73] separately supervised the intermediate layers of a deep network for contour detection. The duality between multiscale contours and hierarchical segmentation [1], [53] was further studied by Maninis et al. [42] by bringing CNNs to the field of generic image segmentation. In this work we explore how to train an FCN for accurately localized dense prediction based on very limited annotation: a single segmented frame.

Semantic Instance Segmentation: Semantic instance segmentation is a relatively new computer vision task which has recently gained increasing attention. In contrast to semantic segmentation or object detection, the goal of instance segmentation is to provide a segmentation mask for each individual instance. The task was first introduced in [19], where they extract both region and foreground features using the R-CNN [16] framework and region proposals. Then, the features are concatenated and classified by an SVM. Several works [8], [9], [75] following that path have been proposed in recent years. There also exist some approaches based on iteration [34], and recurrent neural networks [58]. The recent best-performing methods use fully convolutional position sensitive architectures [36], or a modified Faster-RCNN [57] pipeline, extended to instance segmentation [21]. In contrast to such class-sensitive methods, in which unseen classes are treated as background, our method is class agnostic, and is able to segment

generic objects, given only one annotated example.

Using Semantic Information to Aid Other Computer Vision Tasks: Semantic information is a very relevant cue for the human vision system, and some computer vision algorithms leverage it to aid various tasks. [18] improves reconstruction quality by jointly reasoning about class segmentation and 3D reconstruction. Using a similar philosophy, [38] estimates the depth of each pixel in a scene from a single monocular image guided by semantic segmentation, and improves the results significantly. To the best of our knowledge, we are the first ones to apply instance semantic information to the task of object segmentation in videos.

Conditional Models: Conditional models prove to be a very powerful tool when the feature statistics are complex. In this way, prior knowledge can be introduced by incorporating a dependency to it. [11] builds a conditional random forest to estimate face landmarks whose classifier is dependent on the pose of head. Similarly, [64] proposes to estimate human pose dependent on torso orientation, or human height, which can be a useful cue for the task of pose estimation. The same also applies to boundary detection, [68] proposes to train a series of conditional boundary detectors, and the detectors are weighted differently during test based on the global context of the test image. In this work, we argue that the feature distribution of foreground and background pixels are essentially different, and so a monolithic classifier for the whole image is bound to be suboptimal. Thus, we utilize the conditional classifier to better model the different distributions.

3 ONE-SHOT VIDEO OBJECT SEGMENTATION (OSVOS)

This section describes our algorithm to gradually fine-tune the CNN in order to build a strong appearance model for video object segmentation given the first annotated frame. This was presented in our conference contribution [5]. We will refer to the method as OSVOS, to differentiate it from OSVOS^S (Section 4), in which we use semantic instance segmentation as further guiding signal.

Let us assume that one would like to segment an object in a video, for which the only available piece of information is its foreground/background segmentation in one frame. Intuitively, one could analyze the entity, create a *model*, and search for it in the rest of the frames. For humans, this very limited amount of information is more than enough, and changes in appearance, shape, occlusions, etc. do not pose a significant challenge, because we leverage strong priors: first “It is an object,” and then “It is *this particular* object.” Our method is inspired by this gradual refinement.

We train a Fully Convolutional Neural Network (FCN) for the binary classification task of separating the foreground object from the background. We use two successive training steps: First we train on a large variety of objects, offline, to construct a model that is able to discriminate the general notion of a foreground object, i.e., “It is an object.” Then, at test time, we fine-tune the network for a small number of iterations on the particular instance that we aim to segment, i.e., “It is *this particular* object.” The overview of our method is illustrated in Figure 2.

3.1 End-to-end trainable foreground FCN

Ideally, we would like our CNN architecture to satisfy the following criteria: (i) Accurately localized segmentation output, as

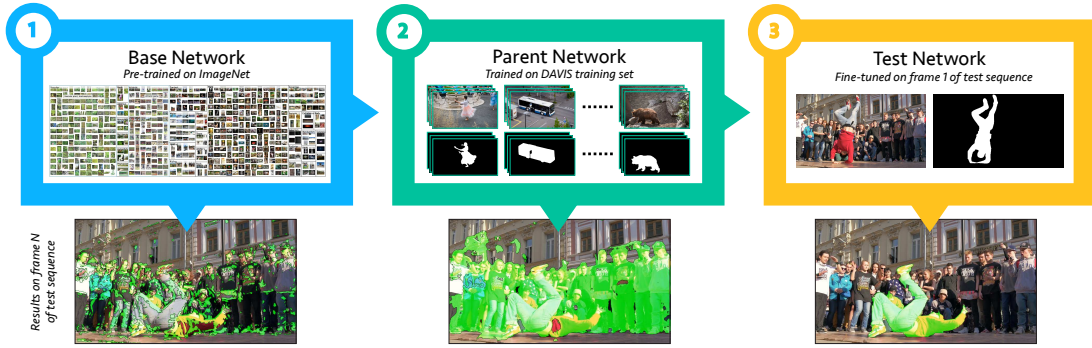


Fig. 2. **Overview of OSVOS:** (1) We start with a pre-trained base CNN for image labeling on ImageNet; its results in terms of segmentation, although conform with some image features, are not useful. (2) We then train a *parent network* on the training set of DAVIS 2016; the segmentation results improve but are not focused on an specific object yet. (3) By fine-tuning on a segmentation example for the specific target object in a single frame, the network rapidly focuses on that target.

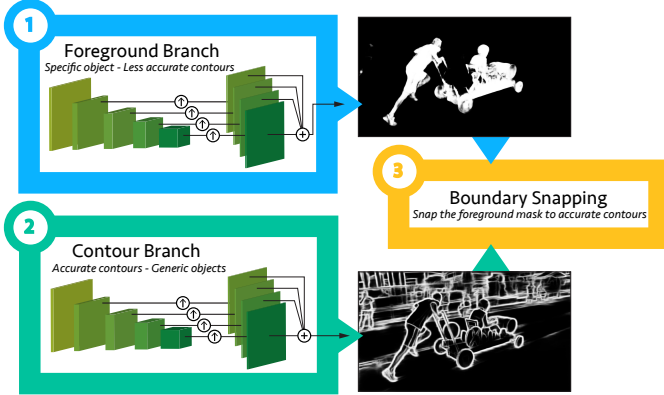


Fig. 3. **Two-stream FCN architecture:** The main foreground branch (1) is complemented by a contour branch (2) which improves the localization of the boundaries (3).

discussed in Section 2, (ii) relatively small number of parameters to train from a limited amount of annotated data, and (iii) relatively fast testing times.

We draw inspiration from the CNN architecture of [43], originally used for biomedical image segmentation. It is based on the VGG [62] network, modified for accurately localized dense prediction (Point i). The fully-connected layers needed for classification are removed (Point ii), and efficient image-to-image inference is performed (Point iii). The VGG architecture consists of groups of convolutional plus Rectified Linear Units (ReLU) [45] layers grouped into 5 stages. Between the stages, pooling operations downscale the feature maps as we go deeper into the network. We connect convolutional layers to form separate skip paths from the last layer of each stage (before pooling). Upscaling operations take place wherever necessary, and feature maps from the separate paths are concatenated to construct a volume with information from different levels of detail. We linearly fuse the feature maps to a single output which has the same dimensions as the image, and we assign a loss function to it. The proposed architecture is shown in Figure 3 (1), foreground branch.

The pixel-wise cross-entropy loss for binary classification (we keep the notation of Xie and Tu [72]) is in this case defined as:

$$\begin{aligned} \mathcal{L}(\mathbf{W}) &= -\sum_j y_j \log P(y_j=1|X; \mathbf{W}) + (1-y_j) \log(1-P(y_j=1|X; \mathbf{W})) \\ &= -\sum_{j \in Y_+} \log P(y_j=1|X; \mathbf{W}) - \sum_{j \in Y_-} \log P(y_j=0|X; \mathbf{W}) \end{aligned}$$

where \mathbf{W} are the standard trainable parameters of a CNN, X is the input image, $y_j \in \{0, 1\}$, $j = 1, \dots, |X|$ is the pixel-wise binary label of X , and Y_+ and Y_- are the positive and negative labeled pixels. $P(\cdot)$ is obtained by applying a sigmoid to the activation of the final layer.

In order to handle the imbalance between the two binary classes, Xie and Tu [72] proposed a modified version of the cost function, originally used for contour detection (we drop \mathbf{W} for the sake of readability):

$$\mathcal{L}_{mod} = -\beta \sum_{j \in Y_+} \log P(y_j=1|X) - (1-\beta) \sum_{j \in Y_-} \log P(y_j=0|X) \quad (1)$$

where $\beta = |Y_-|/|Y|$. Equation 1 allows training for imbalanced binary tasks [31], [42], [43], [72].

3.2 Training details

Offline training: The base CNN of our architecture [62] is pre-trained on ImageNet [60] for image labeling, which has proven to be a very good initialization to other tasks [20], [31], [40], [42], [72], [74]. Without further training, the network is not capable of performing segmentation, as illustrated in Figure 2 (1). We refer to this network as the “*base network*.”

We therefore further train the network on the binary masks of the training set of DAVIS 2016, to learn a generic notion of how to segment objects from their background, their usual shapes, etc. We use Stochastic Gradient Descent (SGD) with momentum 0.9 for 50000 iterations. We augment the data by mirroring and zooming in. The learning rate is set to 10^{-8} , and is gradually decreased. After offline training, the network learns to segment foreground objects from the background, as illustrated in Figure 2 (2). We refer to this network as the “*parent network*.”

Online training/testing: With the parent network available, we can proceed to our main task (“*test network*” in Figure 2): Segmenting a particular entity in a video, given the image and the segmentation of the first frame. We proceed by further training (fine-tuning) the parent network for the particular image/ground-truth pair, and then testing on the entire sequence, using the new

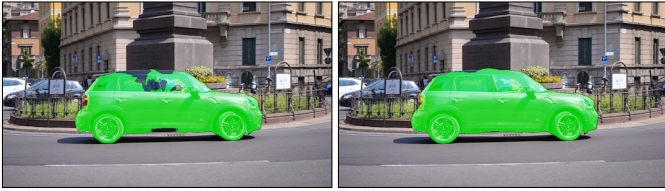


Fig. 4. **Qualitative evolution of the fine tuning:** Results at 10 seconds and 1 minute per sequence.

weights. The timing of our method is therefore affected by two times: the fine-tuning time (once per annotated mask) and the segmentation of all frames (once per frame). In the former we have a trade-off between quality and time: the more iterations we allow the technique to learn, the better results but the longer the user will have to wait for results. The latter does not depend on the training time: OSVOS is able to segment each 480p frame (480×854) in 130 ms.

Regarding the fine-tuning time, we present two different modes: One can either need to fine-tune online, by segmenting a frame and waiting for the results in the entire sequence, or offline, having access to the object to segment beforehand. Especially in the former mode, there is the need to control the amount of time dedicated to training: the more time allocated for fine-tuning, the more the user waits and the better the results are. In order to explore this trade-off, in our experiments we train for a period between 10 seconds and 10 minutes per sequence. Figure 4 shows a qualitative example of the evolution of the results quality depending on the time allowed for fine-tuning. In the experimental evaluation, Figure 12 quantifies this evolution.

Ablation analysis shows that both offline and online training are crucial for good performance: If we perform our online training directly from the base network (ImageNet model), the performance drops significantly. Only dropping the online training for a specific object (using the parent network directly) also yields a significantly worse performance, as already transpired from Figure 2.

3.3 Contour snapping

In the field of image classification [22], [32], [62], where our base network was designed and trained, spatial invariance is a design choice: no matter where an object appears in the image, the classification result should be the same. This is in contrast to the accurate localization of the object contours that we expect in (video) object segmentation. Despite the use of skip connections [20], [40], [43], [73] to minimize the loss of spatial accuracy, we observe that OSVOS' segmentations have some room for improvement in terms of contour localization.

To overcome this limitation, we propose a complementary CNN in a second branch that is trained to detect object contours. Figure 3 shows the global architecture: (1) shows the main foreground branch, where the foreground pixels are estimated; (2) shows the contour branch, which detects all contours in the scene (not only those of the foreground object). This allows us to train offline, without the need to fine-tune on a specific example. We used the exact same architecture in the two branches, but training for different losses. We noticed that jointly training a network with shared layers for both tasks rather degrades the results thus we kept the computations for the two objectives uncorrelated. This allows us to train the contour branch only offline and thus it does not affect the online timing. Since there is need for high recall

in the contours, we train on the PASCAL-Context [44] database, which provides contour annotations for the full scene of an image.

Once we have the estimated object contours, the boundary snapping step (Figure 3 (3)), consists of two different steps:

a) **Superpixel snapping:** It computes superpixels that align to the computed contours (branch 2) by means of an Ultrametric Contour Map (UCM) [1], [53], which we *threshold* at a low strength. We then take a foreground mask (branch 1) and we select superpixels via majority voting (those that overlap with the foreground mask over 50%) to form the final foreground segmentation.

b) **Contour recovery:** It recovers the very thin structures that are lost when snapping to superpixels. It enumerates the connected components of the foreground mask (branch 1), and then matches their contours to the detected contours in branch (2). The connected components whose contour matches the generic contours (branch 2) above a certain tolerance are added to the final result mask.

This refinement process results in a further boost in performance, and is fully modular, meaning that depending on the timing requirements one can choose not to use them, sacrificing accuracy for execution time; since the module comes with a small, yet avoidable computational overhead. Please refer to the timing experiments (Figure 12) for a quantitative evaluation of this trade off: at which range of desired speeds one can afford to use contour snapping.

4 SEMANTIC GUIDANCE (OSVOS^S)

The motivation behind semantic guidance is to improve the model we construct from the first frame with information about the category of the object and the number of instances, e.g. we track two people and a motorbike. We extract the semantic instance information from instance-aware semantic segmentation algorithms. We experiment with three top-performing methods: MNC [9], FCIS [36] and the most recent MaskRCNN [21]. We modify the algorithm and the network architecture to select and propagate the specific instances we are interested in (Section 4.2), and then we adapt the network architecture to include these instance inside the CNN (Section 4.3). The global network overview is first presented in Section 4.1.

4.1 Network Overview

Figure 5 illustrates the structure and workflow of the proposed semantic-aware network. Sharing the common base network (VGG) as the feature extractor, three pixel-wise classifiers are jointly learned.

The first classifier, First-Round Foreground Estimation, is the original OSVOS head, which is purely appearance based, with no knowledge about the semantic segmentation source and produces the first foreground estimation. The result of that classifier and the information from an external semantic instance segmentation system are combined in the semantic selection and propagation steps (Section 4.2) to produce the top matching instances that we refer to as the semantic prior.

The two other classifiers inside the conditional classifier operate on both the features of the common base network and the semantic prior, and are dependent on each other: one is responsible for the pixels with a foreground prior, whereas the other for the background ones. Finally, the two sets of predictions are fused into the final prediction. See Section 4.3.

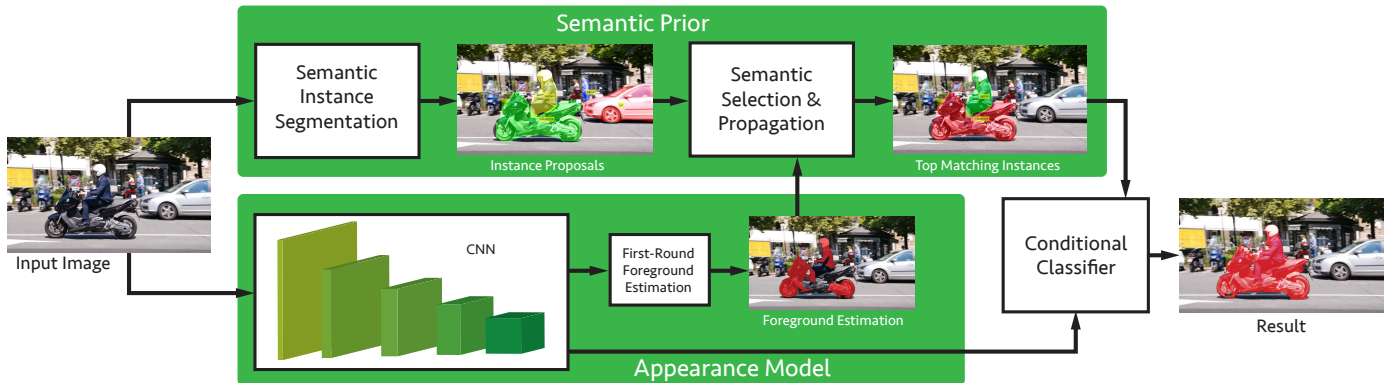


Fig. 5. **Network architecture overview:** Our network is composed of three major components: a base network as the feature extractor, and three classifiers built on top with shared features: a first-round foreground estimator to produce the semantic prior, and two conditional classifiers to model the appearance likelihood.

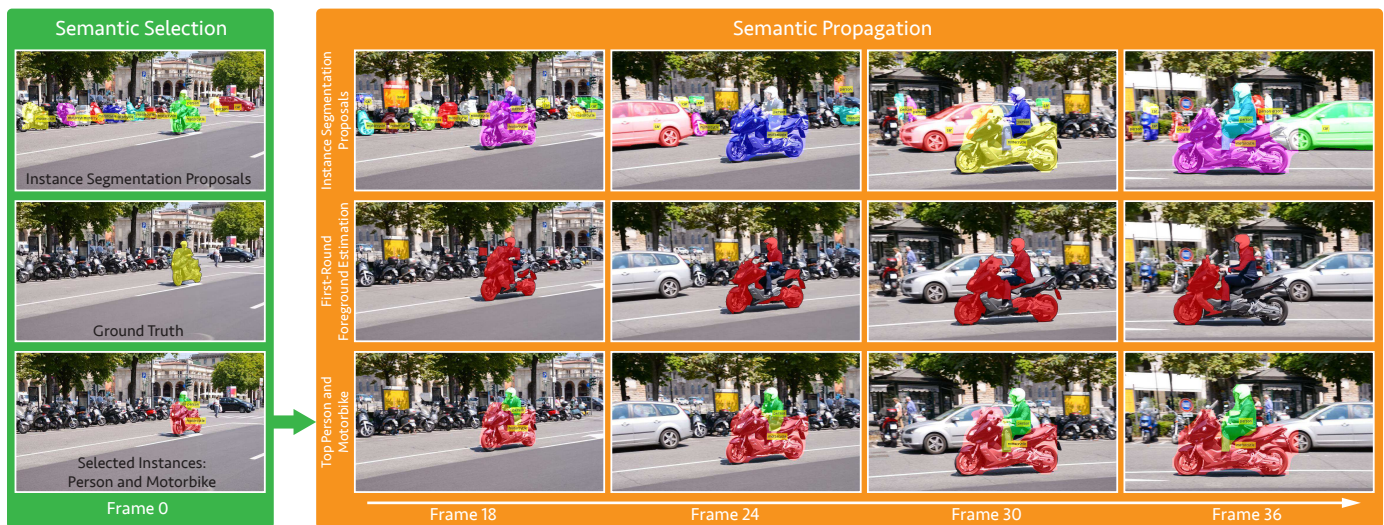


Fig. 6. **Semantic selection and propagation:** Illustrative example of the estimation of the semantics of the object from the first frame (semantic selection) and its propagation to the following frames (semantic propagation).

4.2 Semantic Selection and Semantic Propagation

We leverage a semantic instance segmentation algorithm as an input to estimate the semantics of the object to be segmented. Specifically, we choose MNC [9], FCIS [36], or MaskRCNN [21] as our input instance segmentation algorithms, and we use their publicly available implementations. We show that each of the improvements in instance segmentation is translated in a boost for the task of video object segmentation, which suggests that our method will be able to incorporate future improvements in the field.

The three instance semantic segmentation methods (MNC, FCIS, and MaskRCNN) are multi-stage networks that consist of three major components: shared convolutional layers, region proposal network (RPN), and region-of-interest(ROI)-wise classifiers. We use the available models which are pre-trained on PASCAL for the first one and on COCO for the other two. We note that our method is category agnostic, and the objects to segment do not necessarily need to be part of the PASCAL or COCO category vocabulary, as it will be shown in the experiments.

The output of the instance segmentation algorithm is given as a set of binary masks, together with their category, and their confidence of being a true object. We search for the object of

interest inside the pool of most confident masks: our objective is to find a subset of masks with consistent semantics throughout the video as our semantic prior.

The process can be divided into two stages, namely *semantic selection* and *semantic propagation*. Semantic selection happens in the first frame, where we select the masks that match the object according to the given ground-truth mask (please note that we are in a semi-supervised framework where the true mask of the first frame is given as input). The number of instances and their categories are what we enforce to be consistent throughout the entire video. Figure 6 depicts an example of both steps. Semantic selection, on the left in green, finds that we are interested in a motorbike plus a person (bottom), by overlapping the ground truth (middle) to the instance segmentation proposals (top). There are two cases where semantic selection may fail: a) the objects of interest are not part of the semantic vocabulary of the instance segmenter, and b) the wrong instances are selected by this step. Results show that our classifiers are robust to such failures, preserving high quality outputs in both cases. Thus, a fast greedy search for selecting the instances is sufficient to preserve high performance.

The semantic propagation stage (in orange) occurs at the

following frames, where we propagate the semantic prior we estimated in the first frame to the following ones. No information from future frames is used in this stage. The instance segmentation masks (first row), are filtered using the first-round foreground estimation from the OSVOS head (middle row), and the top matching person and motorbike from the pool are selected (bottom row). In cases that an instance of the selected classes does not overlap with the output of OSVOS, as in cases of occlusions and moving camera, we exclude the particular instance from the semantic prior, for the specific frame.

4.3 Conditional Classifier

Dense labeling using fully convolutional networks is commonly formulated as a per-pixel classification problem. It can be therefore understood as a *global* classifier sliding over the whole image, and assigning either the foreground or background label to each pixel according to a *monolithic* appearance model. In this work, we want to incorporate the semantic prior to the final classification, which will be given as a mask of the most promising instance (or set of instances) in the current frame.

If semantic instance segmentation worked perfectly, we could directly select the best-matching instance to the appearance model, but in reality the results are far from perfect (as we will show in the experiments). We can only, therefore, use the instance segmentation mask as a guidance, or a guess, of what the limits of that instance are, but we still need to perform a refinement step. Our proposed solution to incorporate this mask but still keep the per-pixel classification is to train two classifiers and weigh them according to the confidence we have in that pixel being part of the instance or not. We argue that using a single set of parameters for the whole image is suboptimal.

Formally, for each pixel i , we estimate its probability of being a foreground pixel given the image: $p(i|I)$. The probability can be decomposed into the sum of k conditional probabilities weighted by the prior:

$$p(i|I) = \sum_{k=1}^K p(i|I, k) p(k|I).$$

In our experiments, we use $K = 2$, and we build two conditional classifiers, one focusing on the *instance foreground* pixels, and the other focusing on the *instance background* pixels. The prior term $p(k|I)$ is estimated based on the instance segmentation output. Specifically, a pixel relies more on the *instance foreground* classifier if it is located within the instance segmentation mask; and more importance is given to the *instance background* classifier if it falls out of the instance segmentation mask. In our experiments, we apply a Gaussian filter to spatially smooth the selected masks as our semantic prior.

The conditional classifier is implemented as a layer which can be easily integrated in the network in an end-to-end trainable manner. The layer takes two prediction maps f_1 and f_2 and the weight maps $p(k|I)$ which come from the semantic selection. Without loss of generality, we will assume that $k = 1$ corresponds to the foreground of the semantic prior. For convenience, we set $w = p(k = 1|I)$, and in our case $1 - w = p(k = 2|I)$ (background prior). The inference process is illustrated in Figure 7, where each input element is multiplied by its corresponding weight from the weight map, then summed with the corresponding element in the other map:

$$f_{out}(x, y) = w(x, y) f_1(x, y) + (1 - w(x, y)) f_2(x, y). \quad (2)$$

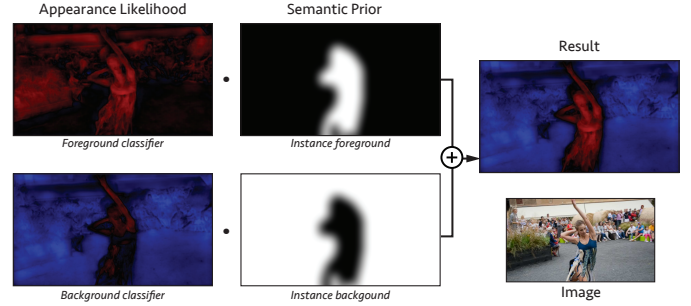


Fig. 7. **Forward pass of the conditional classifier layer:** Red denotes foreground probability, and blue background probability. The output is the weighted sum of the two conditional classifier.

In Equation 2, x and y represent the horizontal and vertical pixel location on a frame. This equation suggest that the decision for the pixels near the selected instances are made by the instance foreground classifier ($f_1(x, y)$), whereas the instance background classifier ($f_2(x, y)$) decides for the rest of the pixels.

Similarly, in the back-propagation step, the gradient from the top g_{top} is propagated to the two parts according to the weight map:

$$g_1(x, y) = w(x, y) g_{top}(x, y)$$

$$g_2(x, y) = (1 - w(x, y)) g_{top}(x, y).$$

The conditional classifier is necessary to incorporate the semantic prior information, in order to make softer decisions. Techniques that can be used as alternatives incorporating only a single classifier, such as masking of the features by the semantic prior, lead to hard decisions guided by the semantics, unable to recover in regions where they are wrong. For example, in Figure 7, the left hand of the dancer is not detected by the semantic prior, and it will be immediately classified as background in the case of feature masking. The background classifier of our proposed method, however, is able to recover the region, correctly classifying it as a foreground.

4.4 Training and Inference

We follow the same ideas as OSVOS to train and test the network, every step enriched with the semantic selection and propagation steps. The parent network is trained using semantic instances that overlap with the ground-truth masks of the DAVIS 2016 training set. Similarly, during online fine-tuning we use the label of the first frame, as well as the outputs of the OSVOS head for the next ones. As was done before, each frame is processed independently of the others. As shown in the experiments, the plug-in of the instance segmentation module dramatically improves the quality of the final segmentation.

5 EXPERIMENTAL VALIDATION

Experimental Setup: We will mainly work on the DAVIS 2016 database [50], using their proposed metrics: region similarity (intersection over union \mathcal{J}), contour accuracy (\mathcal{F} measure), and temporal instability (\mathcal{T}). The dataset contains 50 full-HD annotated video sequences, 30 in the training set and 20 in the validation set. All our results will be trained on the former, evaluated on the latter. As a global comparison metric we will



Fig. 8. **Semantic selection evaluation:** Semantic instances selected by the semantic selection step, with its category overlaid. We observe that in some cases either the semantic labels (h-i) or the number of instances (j) is incorrect. The final results, however, are robust to such mistakes.

	Measure	ImageNet	+OneShot	+Parent	+Semantics	+Superpixels	+Contours					
$\mathcal{J}\&\mathcal{F}$	Mean \mathcal{M} \uparrow	18.9	65.6	<i>46.7</i>	77.8	<i>12.1</i>	86.1	<i>8.4</i>	85.4	<i>0.7</i>	86.5	<i>1.1</i>
	Mean \mathcal{M} \uparrow	17.6	64.6	<i>47.0</i>	77.4	<i>12.8</i>	85.0	<i>7.6</i>	85.5	<i>0.5</i>	85.6	<i>0.1</i>
\mathcal{J}	Recall \mathcal{O} \uparrow	2.3	70.5	<i>68.2</i>	91.0	<i>20.5</i>	96.7	<i>5.7</i>	96.5	<i>0.2</i>	96.8	<i>0.3</i>
	Decay \mathcal{D} \downarrow	1.8	27.8	<i>26.0</i>	17.4	<i>10.4</i>	7.2	<i>10.2</i>	5.9	<i>1.4</i>	5.5	<i>0.3</i>
\mathcal{F}	Mean \mathcal{M} \uparrow	20.3	66.7	<i>46.4</i>	78.1	<i>11.4</i>	87.3	<i>9.2</i>	85.3	<i>2.0</i>	87.5	<i>2.2</i>
	Recall \mathcal{O} \uparrow	2.4	74.4	<i>72.0</i>	92.0	<i>17.6</i>	95.9	<i>3.9</i>	94.1	<i>1.8</i>	95.9	<i>1.8</i>
	Decay \mathcal{D} \downarrow	2.4	26.4	<i>24.0</i>	19.4	<i>7.0</i>	9.3	<i>10.1</i>	6.8	<i>2.5</i>	8.2	<i>1.5</i>
\mathcal{T}	Mean \mathcal{M} \downarrow	46.0	60.9	<i>14.9</i>	33.5	<i>27.4</i>	20.2	<i>13.3</i>	25.1	<i>4.9</i>	21.7	<i>3.4</i>

TABLE 1

Ablation study on DAVIS 2016: From a network pretrained on ImageNet, all improvement steps to the proposed OSVOS^S (right-most column). Numbers in italics show how much the results improve (in blue) or worsen (in red) in that metric with respect to the previous column.

use the mean between \mathcal{J} and \mathcal{F} , as proposed in the DAVIS 2017 challenge [54].

We compare against a large body of very recent semi-supervised state-of-the-art techniques (OnAVOS [70], MSK [49], CTN [28], VPN [26], OFL [67], BVS [46], and FCP [51]) using the pre-computed results provided by the respective authors. For context, we also add the results of the latest unsupervised techniques (ARP [30], FSEG [24], LMP [66], FST [48], NLC [12], MSG [4]).

Moreover, we perform experiments on DAVIS 2017 which contains videos with multiple objects. We compute the results on the *test-dev* set using the submission website provided by the organizers of the challenge. We compare against OnAVOS [70], its submission to the DAVIS 2017 challenge which achieves the fifth place [69] and the other top-performing methods of the challenge [29], [33], [35], [61].

For completeness, we also experiment on the Youtube-objects dataset [25], [55] against those techniques with public segmentation results (OnAVOS [70], OSVOS [5], MSK [49], OFL [67], BVS [46]). We do not take pre-computed evaluation results directly from the paper tables because the benchmarking algorithm is not consistent among the different authors.

Ablation Study: Table 1 shows how much each of the improvements presented builds up to the final result. We start by evaluating the network using only ImageNet pre-trained weights, before including any further training to the pipeline. The results in terms of segmentation ($\mathcal{J}\&\mathcal{F} = 18.9\%$) are completely random (as visually shown in Figure 2). Fine-tuning on the mask of the first frame already boosts the results to competitive levels (+OneShot). By pre-training the parent model, we allow fine-tuning to start from a much more meaningful set of weights, from a problem closer to the final one, so performance increases by

Measure	MNC		FCIS		Mask-RCNN		OSVOS ^S	
	Automatic	Oracle	Automatic	Oracle	Automatic	Oracle		
$\mathcal{J}\&\mathcal{F}$	\mathcal{M} \uparrow	63.7	81.5	73.1	75.1	82.4	82.6	86.5
	\mathcal{M} \uparrow	68.9	81.3	74.3	76.4	82.6	82.8	85.6
	\mathcal{O} \uparrow	85.5	95.8	88.4	92.0	93.7	94.3	96.8
	\mathcal{D} \downarrow	3.3	8.4	1.9	1.8	3.0	2.2	5.5
\mathcal{F}	\mathcal{M} \uparrow	58.5	81.6	71.9	73.7	82.2	82.3	87.5
	\mathcal{O} \uparrow	63.0	93.3	82.8	87.2	88.8	89.5	95.9
	\mathcal{D} \downarrow	3.0	13.6	3.0	3.2	3.4	2.7	8.2
\mathcal{T}	\mathcal{M} \downarrow	30.5	28.4	24.8	23.9	17.4	17.5	21.7

TABLE 2

Semantic propagation: Comparing the automatic selection of instances against an oracle and our final result.

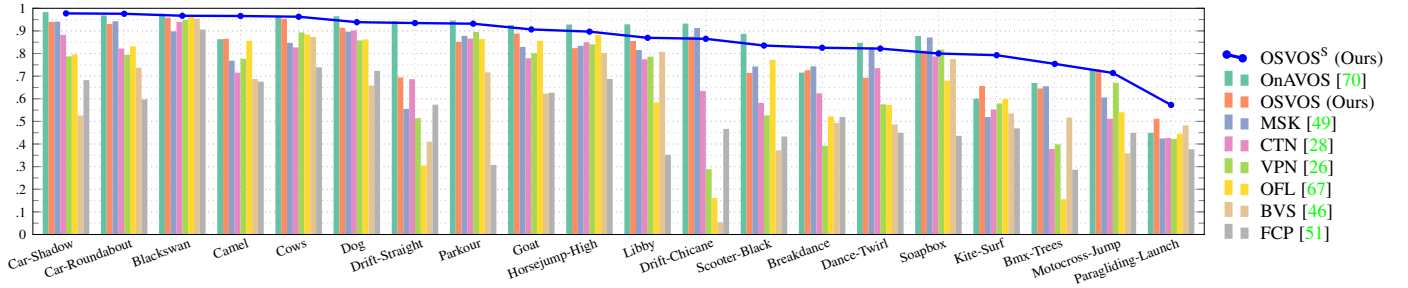
12% (+Parent). Adding semantics and the conditional classifier (+Semantics) plays an important role both in terms of regions and contours ($\mathcal{J}\&\mathcal{F}$), but especially on temporal stability (\mathcal{T}). Snapping to superpixels (+Superpixels) and recovering the contours (+Contours) improve the results around half a point overall, the former especially in terms of \mathcal{J} , the latter in terms of \mathcal{F} , as it stands to reason.

Semantic Selection and Propagation: Figure 8 qualitatively evaluates the semantic-selection algorithm: it displays the selected semantic instances on the first frame of eight videos. Examples (a) and (b) show correct detections in terms of category when a single instance is present. Results (c) to (f) show that the algorithm works also in terms of the quantity of instances when more than one of them is needed. Images (g) to (i) display cases where the category of the object is not present in MS COCO [37] (on which the instance segmentation algorithm was trained), so the *closest semantic* match is used instead. Please note that the precise category is not needed for our algorithm to work, as long as that category is consistent throughout the video (e.g. as long

Measure	Semi-Supervised										Unsupervised					Bound	
	OSVOS ^S	OnAVOS	OSVOS	MSK	CTN	VPN	OFL	BVS	FCP	ARP	FSEG	LMP	NLC	FST	MSG	COB SP	
$\mathcal{J}\&\mathcal{F}$	Mean \mathcal{M} \uparrow	86.5	85.5	80.2	77.5	71.4	67.8	65.7	59.4	53.8	73.4	68.0	67.9	53.7	53.4	52.1	86.8
\mathcal{J}	Mean \mathcal{M} \uparrow	85.6	86.1	79.8	79.7	73.5	70.2	68.0	60.0	58.4	76.2	70.7	70.0	55.1	55.8	53.3	86.5
	Recall \mathcal{O} \uparrow	96.8	96.1	93.6	93.1	87.4	82.3	75.6	66.9	71.5	91.1	83.5	85.0	55.8	64.9	61.6	96.5
\mathcal{F}	Decay \mathcal{D} \downarrow	5.5	5.2	14.9	8.9	15.6	12.4	26.4	28.9	-2.0	7.0	1.5	1.3	12.6	0.0	2.4	2.8
	Mean \mathcal{M} \uparrow	87.5	84.9	80.6	75.4	69.3	65.5	63.4	58.8	49.2	70.6	65.3	65.9	52.3	51.1	50.8	87.1
\mathcal{T}	Recall \mathcal{O} \uparrow	95.9	89.7	92.6	87.1	79.6	69.0	70.4	67.9	49.5	83.5	73.8	79.2	51.9	51.6	60.0	92.4
	Decay \mathcal{D} \downarrow	8.2	5.8	15.0	9.0	12.9	14.4	27.2	21.3	-1.1	7.9	1.8	2.5	11.4	2.9	5.1	2.3
\mathcal{T}	Mean \mathcal{M} \downarrow	21.7	19.0	37.8	21.8	22.0	32.4	22.2	34.7	30.6	39.3	32.8	57.2	42.5	36.6	30.1	27.9
Training Images		2.3k + 83k [§]	87k [§]	2.3k [§]	11k [§]	11.4k [§]	2.3k [§]	1	1	1	—	—	—	—	—	—	—

TABLE 3

DAVIS 2016 Validation: OSVOS^S versus the state of the art (both semi- and un-supervised, and a practical bound). For the number of images, we count those datasets that have some form of segmentation (instance or semantic), and we mark the models pre-trained on Imagenet with [§]. k stands for thousands. The number of images in italics is not directly used to train for the task of video object segmentation, but to train the auxiliary semantic instance segmentation network used by OSVOS^S.

Fig. 9. **DAVIS 2016 Validation:** Per-sequence results of mean region similarity and contour accuracy ($\mathcal{J}\&\mathcal{F}$).

as the camel is always detected as a cow). Last image (j) shows a failure case where two persons are detected when just a single one (albeit upside down) is present, but the algorithm is afterwards robust to this mistake.

Once the semantic selection is done on the first frame, the information is propagated throughout the video. Table 2 quantitatively evaluates this step by comparing our automatic selection of instances against an oracle that selects the best instance in each frame independently. We use three different instance segmentation algorithms (MNC [9], FCIS [7] and MaskRCNN [21]). The results show that in all cases our automatic selection gets very close to the oracle selection (best possible instance), so we are not losing much quality in this step; and this is so in all instance segmentation algorithms, showing that we are robust to the particular algorithm used and so we will be able to incorporate future improvements in this front. The last column shows our final result, which significantly improves the oracle selection, so instance segmentation alone is not enough, as already pointed out in previous sections. For the rest of the paper, we refer to OSVOS^S as the OSVOS^S-MaskRCNN variant of our method.

Comparison to State of the Art in DAVIS 2016: Table 3 shows the comparison of OSVOS^S and OSVOS against a large set of very recent video segmentation algorithms, semi-supervised (using the first segmented frame as input) and unsupervised (only the raw video as input). Apart from the standard metrics of DAVIS 2016 [50], we also add the most recent mean between \mathcal{J} and \mathcal{F} , as used in the 2017 DAVIS Challenge [54].

OSVOS^S is the best performing technique overall, one point above the second semi-supervised technique and 12.6 points above the best unsupervised one. Last column shows the best result one could obtain from picking superpixels from COB [41], [42], a state-of-the-art generic image segmentation algorithm, at a very fine scale. We select the superpixels by snapping the ground-truth masks to them, thus creating a very strong bound. OSVOS^S is only

Attr	OSVOS ^S	OnAVOS	OSVOS	MSK	CTN	VPN	OFL							
LR	89.3	-3.6	89.5	-5.3	80.1	0.1	78.9	-1.8	69.0	3.2	57.7	13.5	45.7	26.7
SV	82.9	6.2	82.3	5.4	74.8	9.1	71.8	9.6	62.5	14.8	58.6	15.3	51.0	24.5
FM	85.2	2.1	84.2	1.9	77.7	3.9	75.0	4.0	65.8	8.7	57.4	16.1	48.7	26.1
CS	89.8	-5.0	88.3	-4.3	80.8	-0.9	76.8	1.1	71.8	-0.6	68.8	-1.5	64.2	2.3
DB	82.8	4.4	75.0	12.4	75.3	5.8	72.5	5.9	60.4	13.0	42.0	30.4	42.8	27.0
MB	82.8	6.8	80.8	8.5	74.7	9.9	72.1	9.9	66.1	9.5	62.1	10.4	53.6	22.0
OCC	86.8	-0.4	84.0	2.1	79.8	0.6	75.8	2.5	70.8	0.8	73.2	-7.7	66.2	-0.7
OV	82.4	5.2	80.8	5.9	71.1	11.4	68.3	11.6	63.9	9.3	53.8	17.5	48.5	21.5

TABLE 4

Attribute-based performance ($\mathcal{J}\&\mathcal{F}$): Impact of the attributes of the sequences on the results. For each attribute, results on the sequences with that particular feature and in italics the gain with respect to those on the set of sequences without the attribute. LR stands for low resolution, SV for scale variation, FM for fast motion, CS for camera shake, DB for dynamic background, MB for motion blur, OCC for occlusions, and OV for object out of view.

0.3 points below the value of this oracle, further highlighting the outstanding quality of our results.

Next, we break down the performance on DAVIS 2016 per sequence. Figure 9 shows the previous state-of-the-art techniques in bars, and OSVOS^S using a line; sorted by the *difficulty* of the sequence for our technique. We see that we outperform the majority of algorithms in the majority of sequences, especially so in the more challenging ones (e.g. Kite-Surf, Bmx-Trees). Please also note that OSVOS^S results are above 70% in all but one sequence and above 80% in all but three, which highlights the robustness of the approach.

Table 4 shows the per-attribute comparison in DAVIS 2016, that is, the mean results on a subset of sequences where a certain challenging attribute is present (e.g. camera shake or occlusions). The increase/decrease of performance when each attribute is not present (small positive/negative numbers in italics) is significantly low, which shows that OSVOS^S is also very robust to the different challenges.

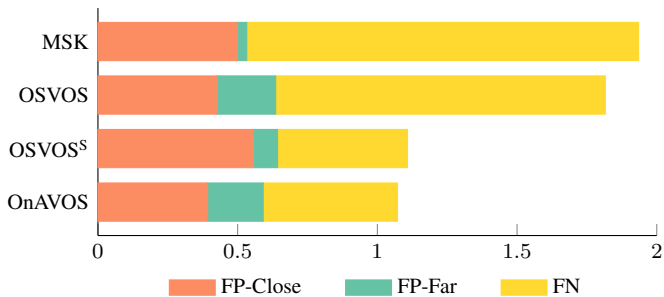


Fig. 10. **Error analysis of our method:** Errors divided into False Positives (FP-Close and FP-Far) and False Negatives (FN). Values are percentage (%) of FP-Close, FP-Far or FN pixels in a sequence.

Number of training images (parent network): To evaluate how many annotated data are needed to retrain a parent network, Table 5 shows the performance of OSVOS^S when using a subset of the DAVIS 2016 training set. We directly used the output of the CNN, without snapping, for efficiency. We randomly selected a fixed percentage of the annotated frames over all videos of the training set, and evaluated using the Region Similarity (\mathcal{J}) metric. We conclude that by using only ~ 200 annotated frames,

Training data	100	200	600	1000	2079
Quality (\mathcal{J})	82.3	84.9	85.2	85.5	85.6

TABLE 5

Amount of training data: Region similarity (\mathcal{J}) as a function of the number of training images for the parent network of OSVOS^S. Full DAVIS 2016 training set is 2079 training data.

we are able to reach almost the same performance than when using the full DAVIS 2016 training split. Thus, we therefore do not require full video annotations for the training procedure, that are often expensive to acquire. Even more, since our method is by definition disregarding temporal information, it is natural that the training data do not require to be temporally coherent.

Misclassified-Pixels Analysis: Figure 10 shows the error analysis of our method. We divide the misclassified pixels in three categories: Close False Positives (FP-Close), Far False Positives (FP-Far) and False Negatives (FN): (i) FP-Close are those near the contour of the object of interest, so contour inaccuracies, (ii) FP-Far reveal if the method detects other objects or blobs apart from the object of interest, and (iii) FN tell us if we miss a part of the object during the sequence. The measure in the plot is the percentage of pixels in a sequence that fall in a certain category.

The main strength of OSVOS^S compared to OSVOS and MSK is considerably reducing the number of false negatives. We believe this is due to OSVOS^S's ability to *complete* the object of interest when parts that were occluded in the first frame become visible, thanks to the semantic concept of instance. On the other hand, the output of the instance segmentation network that we are currently using, FCIS [7], is not very precise on the boundaries of the objects, and even though our conditional classifier is able to recover in part, FP-Close is slightly worse than that of the competition. On the plus side, since the instance segmentation is an independent input to our algorithm, we will probably directly benefit from better instance segmentation algorithms.

Performance Decay: As indicated by the \mathcal{J} -Decay and \mathcal{F} -Decay values in Table 3, OSVOS^S exhibits a better ability than

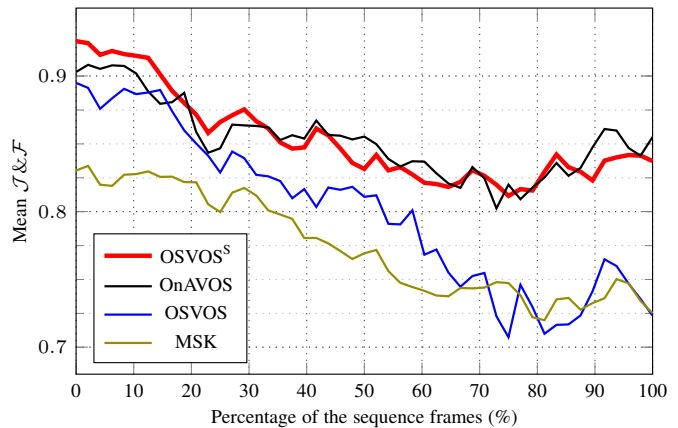


Fig. 11. **Decay of the quality with time:** Performance of various methods with respect to the *time* axis.

OSVOS and MSK to maintain performance as frames evolve, and we interpret that this is so thanks to the injected semantic prior. The performance decay is similar to that of OnAVOS, even though it performs a costly iterative algorithm which fine-tunes the result to various frames of the sequence. Our method, on the other hand, uses the information of the first frame only, and keeps the quality throughout the sequence.

To further highlight this result and analyze it more in detail, Figure 11 shows the evolution of \mathcal{J} as the sequence advances, to examine how the performance drops over time. Since the videos in DAVIS 2016 are of different length, we normalize them to $[0, 100]$ as a percentage of the sequence length. We then compute the mean \mathcal{J} curve among all video sequences. As it can be seen from Figure 11, our method is significantly more stable in terms of performance drop compared to OSVOS and MSK, and has a similar curve than OnAVOS.

We also report the lowest point of the curve which indicates the worst performance across the video. Based on this metrics, OSVOS^S is at 82.0, while for semantic-blind methods, the numbers are 81.0, 73.7, and 69.8.

The results therefore confirm that the semantic prior we introduce can mitigate the performance drop caused by appearance change, while maintaining high fidelity in details. The semantic information is particularly helpful in the later stage of videos where dramatic appearance changes with respect to the first frame are more probable.

Speed: The computational efficiency of video object segmentation is crucial for the algorithms to be usable in practice. OSVOS^S can adapt to different timing requirements, providing progressively better results the more time we can afford, by letting the fine-tuning algorithm at test time do more or fewer iterations. As introduced before, OSVOS^S's time can be divided into the fine-tuning time plus the time to process each frame independently.

To compare to other techniques, we will evaluate the mean computing time per frame: fine-tuning time (done once per sequence) averaged over the length of that sequence, plus the forward pass on each frame.

Figure 12 shows the quality of the result with respect to the time it takes to process each 480p frame. The computation time for our method has been obtained using an NVidia Titan X GPU and for other methods the timing reported in their publications has been used. Our techniques are represented by curves: OSVOS^S

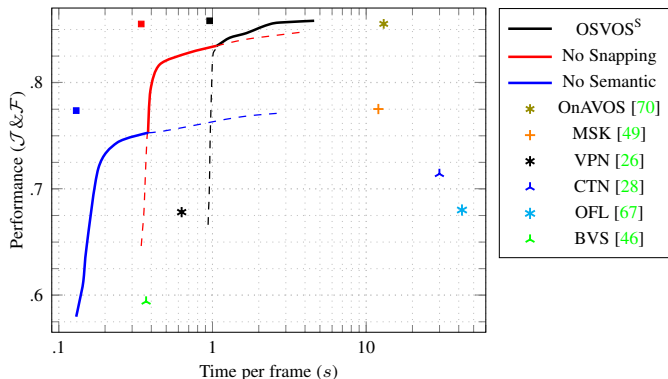


Fig. 12. **Quality versus timing:** $\mathcal{J}\&\mathcal{F}$ with respect to the processing time per frame.

(—), without boundary snapping (—), and without semantics (—), which show the gain in quality with respect to the fine-tuning time. The best results come at the price of adding the semantics or the snapping cost, so depending on the needed speed, one of the three modes can be chosen. Dashed lines represent the regimes of each technique that are not in the Pareto front, i.e. where it is better to choose another mode within our techniques (faster for the same quality or best quality for the same speed).

Since OSVOS^S processes frames independently, one could also perform the fine-tuning offline, by training on a picture of the object to be segmented beforehand (e.g. take a picture of a sports player before a match). In this scenario, OSVOS^S can process each frame by one forward pass of the CNN (▪ | ▪ | ▪), and so be considerably fast.

Compared to other techniques, our techniques are faster and/or more accurate at all regimes, from fast modes: 75.1 versus 59.4 of BVS (▲) at 300 milliseconds, to high-quality regimes: same performance than OnAVOS (*) but an order of magnitude faster (2.5 versus 12 seconds). The trade-off between performance and speed in video object segmentation has been largely ignored (or purposely hidden) in the literature although we believe it is of critical importance, and so we encourage future research to evaluate their methods in this performance-vs-speed plane.

Comparison to State of the Art in Youtube-Objects: For completeness, we also do experiments on Youtube-objects [25], [55], without changing any parameter from our algorithm nor retraining the parent network. Table 6 shows the results of the quantitative evaluation against the rest of techniques. OSVOS^S obtains the best results overall, being two points better than the runner up; and having the best results in eight out of ten categories. These experiments show the robustness and generality of our approach even to domain (dataset) shifts.

Multi-object video segmentation in DAVIS 2017: We test OSVOS^S in the more challenging DAVIS 2017 dataset where multiple objects have to be segmented in the same video sequence. We apply our method as is, pre-processing every object in a sequence independently. Table 7 illustrates the results obtained in the test-dev set of the dataset, compared to the top-performing methods of the DAVIS challenge, and to our direct competitor (OnAVOS). Even though our method is not specifically designed to handle multiple object instances, we achieve competitive results (comparable to the third entry), and we outperform OnAVOS. Our method falls behind the two first entries as it is not optimized to segment multiple objects, and is uses a single model, without the bells and whistles that naturally come with challenge submissions.

Category	OSVOS ^S	OnAVOS	OSVOS	MSK	OFL	BVS
Aeroplane	90.4	87.7	88.2	86.0	89.9	86.8
Bird	87.0	85.7	85.7	85.6	84.2	80.9
Boat	83.6	78.5	77.5	78.8	74.0	65.1
Car	87.9	86.1	79.6	78.8	80.9	68.3
Cat	80.7	80.5	70.8	70.1	68.3	55.8
Cow	79.3	77.9	77.8	77.7	79.8	69.9
Dog	82.5	80.8	81.3	79.2	76.6	68.0
Horse	73.9	72.1	72.8	71.7	72.6	58.9
Motorbike	79.3	72.0	73.5	65.6	73.7	60.5
Train	87.1	84.0	75.7	83.5	76.3	65.2
Mean	83.2	80.5	78.3	77.7	77.6	67.9

TABLE 6
Youtube-Objects evaluation: Per-category and overall mean intersection over union (\mathcal{J}).

Method	Test-Dev $\mathcal{J}\&\mathcal{F}$
Apata [29]	66.6
Lixx [35]	66.1
Wangzhe [33]	57.7
Lalafine123 [61]	57.4
Voiglaender [69]	56.5
OnAVOS [70]	52.8
OSVOS ^S	57.5

TABLE 7
DAVIS 2017 evaluation: Performance of OSVOS^S compared to the DAVIS 2017 challenge winners, on the test-dev set. Our single model achieves competitive results.

Instance segmentation quality: In this section we analyze the influence of the quality of the instance segmentation method in our final result. To this end, we use three different methods, i.e. MNC [9], FCIS [36], and Mask-RCNN [21]. Developments to the field over the last two years have lead to competitive results on COCO [37] test-dev, with resulting Average Precision (AP) varying from 24.6% for MNC, to 33.6% and 37.1% for FCIS and Mask-RCNN, respectively. Table 8 shows the performance gains obtained by using a different instance segmentation method within the OSVOS^S pipeline in three different datasets. Results suggest that our method is able to incorporate improved instance segmentation results, and directly translates them into more accurate results for video object segmentation. The improvements are particularly large for DAVIS 2017, where there is still room for improvement.

Qualitative Results: Figure 13 and Figure 14 show some qualitative results of OSVOS^S in DAVIS 2016 and Youtube-Objects, respectively. The first column shows the ground-truth mask used as input to our algorithm (in red). The rest of the columns show our segmented results in the following frames. These visual results qualitatively corroborate the robustness of our approach to occlusions, dynamic background, change of appearance, etc.

Limitations of OSVOS^S: Both OSVOS and OSVOS^S are very practical for applications due to their accuracy, and their frame-independent design which comes with increased speed compared to competing methods. Limitations of OSVOS mainly regard appearance of objects, such as similar objects, dynamic changes in appearance and viewpoint, and are successfully tackled by introducing the coarse instance segmentation input in OSVOS^S. False positives can be successfully tackled by introducing optical



Fig. 13. **Qualitative results on DAVIS 2016:** OSVOS^s results on a variety of representatives sequences. The input to our algorithm is the ground truth of the first frame (red). Outputs of all frames (green) are produced independent of each other.

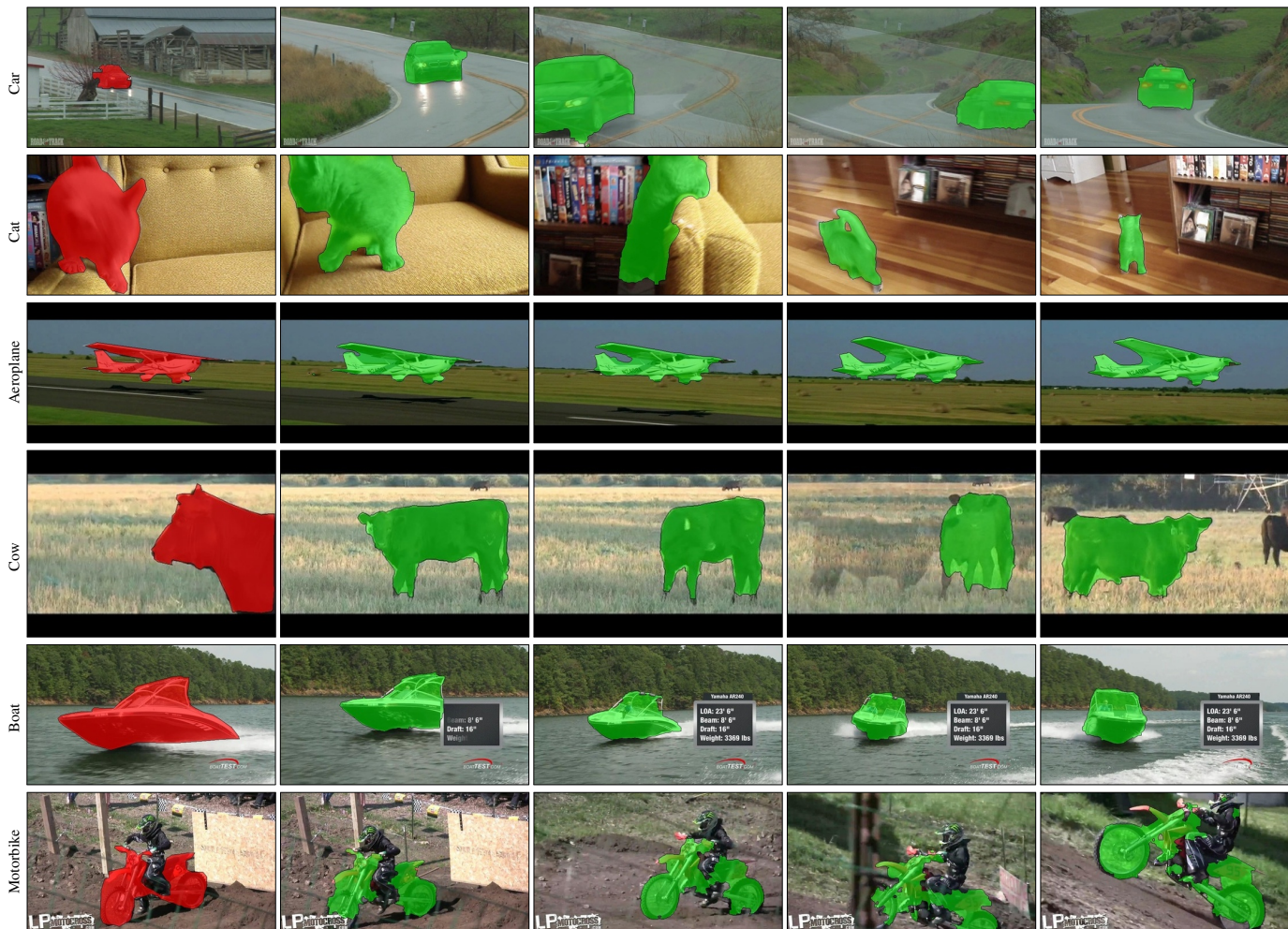


Fig. 14. **Qualitative results on Youtube-Objects:** OSVOS^S results on a variety of representatives sequences. The input to our algorithm is the ground truth of the first frame (red). Outputs of all frames (green) are produced independent of each other.

Dataset	OSVOS ^S			
	<i>Mask-RCNN</i>	<i>FCIS</i>	<i>MNC</i>	OSVOS
DAVIS 2016	86.5	86.0	83.5	80.2
Youtube-Objects	83.2	82.5	80.8	78.3
DAVIS 2017	57.5	53.7	51.5	48.7

TABLE 8

Performance vs. instance segmentation quality: Evaluation with respect to the instance segmentation algorithm.

flow models [29], whereas [35] handle false negatives by introducing a re-identification module, with the cost of extra processing time. Limitations regarding out-of-vocabulary instances are handled well by our method, however, that may not transfer to other domains with uncommon objects or parts of objects.

6 CONCLUSIONS

This paper presents Semantic One-Shot Video Object Segmentation (OSVOS^S), a semi-supervised video object segmentation technique that processes each frame independently and thus ignores the temporal information and redundancy of a video sequence. This has the inherent advantage of being robust to object occlusions, lost frames, etc, while keeping execution speed low.

OSVOS^S shows state-of-the-art results in both DAVIS 2016 and Youtube-Objects at the whole range of operating speeds. It is significantly faster and/or better performing than the competition: 75.1 versus 59.4 at 300 milliseconds per frame, or 4.5 versus 12 seconds at the best performance (86.5 vs 85.5).

To do so, we build a powerful appearance model of the object from a single segmented frame. In contrast to most deep learning approaches, that often require a huge amount of training data in order to solve a specific problem, and in line with humans, that can solve similar challenges with only a single training example; we demonstrate that OSVOS^S can reproduce this capacity of one-shot learning in a machine: Based on a parent network architecture pre-trained on a generic video segmentation dataset, we fine-tune it on merely one training sample.

OSVOS^S also leverages an instance segmentation algorithm that provides a semantic prior to guide the appearance model computed on the first frame. This adds robustness to appearance changes of the object and in practice helps in keeping the quality throughout a longer period of the video.

The appearance model is combined with the semantic prior by means of a new conditional classifier as a trainable module in a CNN.

ACKNOWLEDGMENTS

Research funded by the EU Framework Programme for Research and Innovation Horizon 2020 (Grant No. 645331, EurEyeCase), and by the Swiss Commission for Technology and Innovation (CTI, Grant No. 19015.1 PFES-ES, NeGeVA). The authors gratefully acknowledge support by armasuisse and thank NVidia Corporation for donating the GPUs used in this project.

REFERENCES

- [1] P. Arbeláez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *TPAMI*, 33(5):898–916, 2011. 3, 5
- [2] G. Bertasius, J. Shi, and L. Torresani. High-for-low and low-for-high: Efficient boundary detection from deep object features and its applications to high-level vision. In *ICCV*, 2015. 2
- [3] G. Bertasius, J. Shi, and L. Torresani. Semantic segmentation with boundary neural fields. In *CVPR*, 2016. 2
- [4] T. Brox and J. Malik. Object segmentation by long term analysis of point trajectories. In *ECCV*, 2010. 1, 8
- [5] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. One-shot video object segmentation. In *CVPR*, 2017. 3, 8
- [6] J. Chang, D. Wei, and J. W. Fisher III. A video representation using temporal superpixels. In *CVPR*, 2013. 2, 3
- [7] J. Dai, K. He, Y. Li, S. Ren, and J. Sun. Instance-sensitive fully convolutional networks. In *ECCV*, 2016. 2, 3, 9, 10
- [8] J. Dai, K. He, and J. Sun. Convolutional feature masking for joint object and stuff segmentation. In *CVPR*, 2015. 3
- [9] J. Dai, K. He, and J. Sun. Instance-aware semantic segmentation via multi-task network cascades. In *CVPR*, 2016. 2, 3, 5, 6, 9, 11
- [10] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: Object detection via region-based fully convolutional networks. In *NIPS*, 2016. 2
- [11] M. Dantone, J. Gall, G. Fanelli, and L. Van Gool. Real-time facial feature detection using conditional regression forests. In *CVPR*, 2012. 3
- [12] A. Faktor and M. Irani. Video segmentation by non-local consensus voting. In *BMVC*, 2014. 8
- [13] Q. Fan, F. Zhong, D. Lischinski, D. Cohen-Or, and B. Chen. Jumpcut: Non-successive mask transfer and interpolation for video cutout. *ACM Trans. Graph.*, 34(6), 2015. 2, 3
- [14] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *TPAMI*, 35(8):1915–1929, 2013. 2, 3
- [15] R. Girshick. Fast R-CNN. In *ICCV*, 2015. 1, 3
- [16] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 1, 3
- [17] M. Grundmann, V. Kwatra, M. Han, and I. A. Essa. Efficient hierarchical graph-based video segmentation. In *CVPR*, 2010. 1, 2, 3
- [18] C. Hane, C. Zach, A. Cohen, R. Angst, and M. Pollefeys. Joint 3d scene reconstruction and class segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 97–104, 2013. 3
- [19] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In *ECCV*, 2014. 3
- [20] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *CVPR*, 2015. 2, 4, 5
- [21] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *ICCV*, 2017. 2, 3, 5, 6, 9, 11
- [22] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 3, 5
- [23] H. V. Hopwood. *Living Pictures: Their History, Photo-Production and Practical Working*. Optician & Photographic Trades Review, 1899. 1
- [24] S. Jain, B. Xiong, and K. Grauman. Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. In *CVPR*, 2017. 8
- [25] S. D. Jain and K. Grauman. Supervoxel-consistent foreground propagation in video. In *ECCV*, 2014. 2, 8, 11
- [26] V. Jampani, R. Gadda, and P. V. Gehler. Video propagation networks. In *CVPR*, 2017. 2, 8, 9, 11
- [27] W.-D. Jang and C.-S. Kim. Online video object segmentation via convolutional trident network. In *CVPR*, 2017. 1
- [28] W.-D. Jang and C.-S. Kim. Online video object segmentation via convolutional trident network. In *CVPR*, 2017. 2, 8, 9, 11
- [29] A. Khoreva, R. Benenson, E. Ilg, T. Brox, and B. Schiele. Lucid data dreaming for object tracking. *The 2017 DAVIS Challenge on Video Object Segmentation - CVPR Workshops*, 2017. 8, 11, 13
- [30] Y. Koh and C.-S. Kim. Primary object segmentation in videos based on region augmentation and reduction. In *CVPR*, 2017. 8
- [31] I. Kokkinos. Pushing the boundaries of boundary detection using deep learning. In *ICLR*, 2016. 2, 4
- [32] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1, 3, 5
- [33] T.-N. Le, K.-T. Nguyen, M.-H. Nguyen-Phan, T.-V. Ton, T.-A. N. (2), X.-S. Trinh, Q.-H. Dinh, V.-T. Nguyen, A.-D. Duong, A. Sugimoto, T. V. Nguyen, and M.-T. Tran. Instance re-identification flow for video object segmentation. *The 2017 DAVIS Challenge on Video Object Segmentation - CVPR Workshops*, 2017. 8, 11
- [34] K. Li, B. Hariharan, and J. Malik. Iterative instance segmentation. In *CVPR*, pages 3659–3667, 2016. 3
- [35] X. Li, Y. Qi, Z. Wang, K. Chen, Z. Liu, J. Shi, P. Luo, C. C. Loy, and X. Tang. Video object segmentation with re-identification. *The 2017 DAVIS Challenge on Video Object Segmentation - CVPR Workshops*, 2017. 8, 11, 13
- [36] Y. Li, H. Qi, J. Dai, X. Ji, and Y. Wei. Fully convolutional instance-aware semantic segmentation. In *CVPR*, 2017. 2, 3, 5, 6, 11
- [37] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Zitnick. Microsoft COCO: Common Objects in Context. In *ECCV*, 2014. 8, 11
- [38] B. Liu, S. Gould, and D. Koller. Single image depth estimation from predicted semantic labels. In *CVPR*, 2010. 3
- [39] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, and S. Reed. SSD: Single shot multibox detector. In *ECCV*, 2016. 1
- [40] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2, 3, 4, 5
- [41] K. Maninis, J. Pont-Tuset, P. Arbeláez, and L. V. Gool. Convolutional oriented boundaries: From image segmentation to high-level tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2017. 2, 3, 9
- [42] K. Maninis, J. Pont-Tuset, P. Arbeláez, and L. Van Gool. Convolutional oriented boundaries. In *ECCV*, 2016. 3, 4, 9
- [43] K. Maninis, J. Pont-Tuset, P. Arbeláez, and L. Van Gool. Deep retinal image understanding. In *MICCAI*, 2016. 2, 4, 5
- [44] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, 2014. 5
- [45] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, pages 807–814, 2010. 4
- [46] N. Nicolas Märki, F. Perazzi, O. Wang, and A. Sorkine-Hornung. Bilateral space video segmentation. In *CVPR*, 2016. 2, 3, 8, 9, 11
- [47] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *ICCV*, 2015. 2, 3
- [48] A. Papazoglou and V. Ferrari. Fast object segmentation in unconstrained video. In *ICCV*, 2013. 8
- [49] F. Perazzi, A. Khoreva, R. Benenson, B. Schiele, and A. Sorkine-Hornung. Learning video object segmentation from static images. In *CVPR*, 2017. 1, 2, 3, 8, 9, 11
- [50] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016. 2, 7, 9
- [51] F. Perazzi, O. Wang, M. Gross, and A. Sorkine-Hornung. Fully connected object proposals for video segmentation. In *JCCV*, 2015. 2, 3, 8, 9
- [52] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollár. Learning to refine object segments. In *ECCV*, 2016. 3
- [53] J. Pont-Tuset, P. Arbeláez, J. T. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping for image segmentation and object proposal generation. *TPAMI*, 2017. 3, 5
- [54] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-Hornung, and L. Van Gool. The 2017 DAVIS challenge on video object segmentation. *arXiv:1704.00675*, 2017. 2, 8, 9
- [55] A. Prest, C. Leistner, J. Civera, C. Schmid, and V. Ferrari. Learning object class detectors from weakly annotated video. In *CVPR*, 2012. 2, 8, 11
- [56] S. A. Ramakanth and R. V. Babu. Seamseg: Video object segmentation using patch seams. In *CVPR*, 2014. 1, 2, 3
- [57] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 3
- [58] B. Romera-Paredes and P. H. S. Torr. Recurrent instance segmentation. In *ECCV*, 2016. 3
- [59] O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 3
- [60] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 2, 3, 4
- [61] A. Shaban, A. Firl, A. Humayun, J. Yuan, X. Wang, P. Lei, N. Dhandha, B. Boots, J. M. Rehg, and F. Li. Multiple-instance video segmentation with sequence-specific object proposals. *The 2017 DAVIS Challenge on Video Object Segmentation - CVPR Workshops*, 2017. 8, 11
- [62] K. Simonyan and A. Zisserman. Very deep convolutional networks for

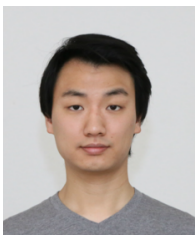
- large-scale image recognition. In *ICLR*, 2015. 1, 3, 4, 5
- [63] S. Stampfer. *Die stroboscopischen Scheiben; oder, Optischen Zauber-scheiben: Deren Theorie und wissenschaftliche Anwendung*. 1833. 1
- [64] M. Sun, P. Kohli, and J. Shotton. Conditional regression forests for human pose estimation. In *CVPR*, 2012. 3
- [65] A. M. Tekalp. *Digital video processing*. Prentice Hall Press, 2015. 1
- [66] P. Tokmakov, K. Alahari, and C. Schmid. Learning motion patterns in videos. In *CVPR*, 2017. 8
- [67] Y.-H. Tsai, M.-H. Yang, and M. J. Black. Video segmentation via object flow. In *CVPR*, 2016. 1, 8, 9, 11
- [68] J. R. Uijlings and V. Ferrari. Situational object boundary detection. In *CVPR*, 2015. 3
- [69] P. Voigtlaender and B. Leibe. Online adaptation of convolutional neural networks for the 2017 davis challenge on video object segmentation. *The 2017 DAVIS Challenge on Video Object Segmentation - CVPR Workshops*, 2017. 8, 11
- [70] P. Voigtlaender and B. Leibe. Online adaptation of convolutional neural networks for video object segmentation. In *BMVC*, 2017. 2, 8, 9, 11
- [71] Y. Wang, Y.-q. Zhang, and J. Ostermann. *Video Processing and Communications*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 1st edition, 2001. 1
- [72] S. Xie and Z. Tu. Holistically-nested edge detection. In *ICCV*, 2015. 2, 4
- [73] S. Xie and Z. Tu. Holistically-nested edge detection. *International Journal of Computer Vision*, pages 1–16, 2017. 3, 5
- [74] J. Yang, B. Price, S. Cohen, H. Lee, and M.-H. Yang. Object contour detection with a fully convolutional encoder-decoder network. In *CVPR*, 2016. 2, 3, 4
- [75] S. Zagoruyko, A. Lerer, T.-Y. Lin, P. O. Pinheiro, S. Gross, S. Chintala, and P. Dollár. A multipath network for object detection. *arXiv preprint arXiv:1604.02135*, 2016. 3



Kevis-Kokitsi Maninis is a PhD candidate at ETHZ, Switzerland, in Prof. Luc Van Gool's Computer Vision Lab (2015). He received the Diploma degree in Electrical and Computer Engineering from National Technical University of Athens (NTUA) in 2014. He worked as undergraduate research assistant in the Signal Processing and Computer Vision group of NTUA (2013-2014).



Sergi Caelles is a Ph.D. candidate at ETHZ, Switzerland, in Prof. Luc Van Gool's Computer Vision Lab (2016). He received the degree in Electrical Engineering and the M.Sc. in Telecommunications Engineering from the Universitat Politècnica de Catalunya, BarcelonaTech (UPC). He worked at Bell Laboratories, New Jersey (USA) in 2014. His research interest include computer vision with special focus on video object segmentation and deep learning.



Yuhua Chen is a PhD candidate at ETHZ, Switzerland, in Prof. Luc Van Gool's Computer Vision Lab (2015). He received a B.Sc in Physics from the University of Science and Technology of China (USTC) in 2013, and M.Sc in Electrical Engineering and Information Technology from ETH Zürich in 2015. His research interests lie in deep learning for semantic segmentation and object detection.



Jordi Pont-Tuset is a post-doctoral researcher at ETHZ, Switzerland, in Prof. Luc Van Gool's Computer Vision Lab (2015). He received the degree in Mathematics in 2008, the degree in Electrical Engineering in 2008, the M.Sc. in Research on Information and Communication Technologies in 2010, and the Ph.D with honors in 2014; all from the Universitat Politècnica de Catalunya, BarcelonaTech (UPC). He worked at Disney Research, Zürich (2014).



Laura Leal-Taixé is leading the Dynamic Vision and Learning group at the Technical University of Munich, Germany. She received her Bachelor and Master degrees in Telecommunications Engineering from the Technical University of Catalonia (UPC), Barcelona. She did her Master Thesis at Northeastern University, Boston, USA and received her PhD degree (Dr.-Ing.) from the Leibniz University Hannover, Germany. During her PhD she did a one-year visit at the Vision Lab at the University of Michigan, USA. She also spent two years as a postdoc at the Institute of Geodesy and Photogrammetry of ETH Zurich, Switzerland and one year at the Technical University of Munich. Her research interests are dynamic scene understanding, in particular multiple object tracking and segmentation, as well as machine learning for video analysis.



Daniel Cremers received Bachelor degrees in Mathematics (1994) and Physics (1994), and a Master's degree in Theoretical Physics (1997) from the University of Heidelberg. In 2002 he obtained a PhD in Computer Science from the University of Mannheim, Germany. Since 2009 he holds the chair for Computer Vision and Pattern Recognition at the Technical University, Munich. His publications received several awards and he has obtained numerous and prestigious funding grants. He has served as area chair (associate editor) for ICCV, ECCV, CVPR, ACCV, IROS, etc, and as program chair for ACCV 2014. He serves as general chair for the European Conference on Computer Vision 2018 in Munich. On March 1st 2016, Prof. Cremers received the Leibniz Award 2016, the biggest award in German academia.



Luc Van Gool got a degree in electromechanical engineering at the Katholieke Universiteit Leuven in 1981. Currently, he is professor at the Katholieke Universiteit Leuven, Belgium, and the ETHZ, Switzerland. He leads computer vision research at both places. His main interests include 3D reconstruction and modeling, object recognition, and tracking, and currently especially their confluence in the creation of autonomous cars. On the latter subject, he leads a large-scale project funded by Toyota. He has authored over 300 papers in this field. He has been a program chair or general chair of several major computer vision conferences. He received several Best Paper awards, incl. a David Marr prize. In 2015, he received the 5-yearly excellence prize of the Flemish Fund for Scientific Research and, in 2016, a Koenderink Award. In 2017 he was nominated one of the main Tech Pioneers in Belgium by business journal 'De Tijd', and one of the 100 Digital Shapers of 2017 by Digitalswitzerland. He is a co-founder of 10 spin-off companies.