

Unseen Object Segmentation in Videos via Transferable Representations

Yi-Wen Chen^{1,2}, Yi-Hsuan Tsai³, Chu-Ya Yang¹,
Yen-Yu Lin¹, Ming-Hsuan Yang^{4,5}

¹Academia Sinica, ²National Taiwan University, ³NEC Laboratories America,
⁴University of California, Merced, ⁵Google Cloud

Abstract. In order to learn object segmentation models in videos, conventional methods require a large amount of pixel-wise ground truth annotations. However, collecting such supervised data is time-consuming and labor-intensive. In this paper, we exploit existing annotations in source images and transfer such visual information to segment videos with unseen object categories. Without using any annotations in the target video, we propose a method to jointly mine useful segments and learn feature representations that better adapt to the target frames. The entire process is decomposed into two tasks: 1) solving a submodular function for selecting object-like segments, and 2) learning a CNN model with a transferable module for adapting seen categories in the source domain to the unseen target video. We present an iterative update scheme between two tasks to self-learn the final solution for object segmentation. Experimental results on numerous benchmark datasets show that the proposed method performs favorably against the state-of-the-art algorithms.

1 Introduction

Nowadays, video data can be easily accessed and visual analytics has become an important task in computer vision. In this line of research, video object segmentation is one of the effective approaches to understand visual contents that facilitates various applications, such as video editing, content retrieval, and object identification. While conventional methods rely on the supervised learning strategy to effectively localize and segment objects in videos, collecting such ground truth annotations is expensive and cannot scale well to a large amount of videos.

Recently, weakly-supervised methods for video object segmentation [40,42,31,41] have been developed to relax the need for annotations where only class-level labels are required. These approaches have significantly reduced the labor-intensive step of collecting pixel-wise training data on target categories. However, these categories are pre-defined and thus the trained model cannot be directly applied to unseen categories in other videos, and annotating additional categories would require more efforts, which is not scalable in practice. In this paper, we propose an algorithm to reduce efforts in both annotating pixel-level and class-level ground truths for unseen categories in videos.

To this end, we make use of existing pixel-level annotations in images from the PASCAL VOC dataset [4] with pre-defined categories, and design a framework to transfer this knowledge to unseen videos. That is, the proposed method is able to learn useful representations for segmentation from the data in the image domain and adapt these representations to segment objects in videos regardless of whether their categories are covered in the PASCAL VOC dataset. Thus, while performing video object segmentation, our algorithm does not require annotations in any forms, such as pixel-level or class-level ground truths.

We formulate the object segmentation problem for unseen categories as a joint objective of mining useful segments from videos while learning transferable knowledge from image representations. Since annotations are not provided in videos, we design an energy function to discover object-like segments in videos based on the feature representations learned from the image data. We then utilize these discovered segments to refine feature representations in a convolutional neural network (CNN) model, while a transferable module is developed to learn the relationships between multiple seen categories in images and the unseen category in videos. By jointly considering both energy functions for mining better segments while learning transferable representations, we develop an iterative optimization method to self-guide the video object segmentation process. We also note that the proposed framework is flexible as we can input either weakly-labeled or unlabeled videos.

To validate the proposed method, we conduct experiments on benchmark datasets for video object segmentation. First, we evaluate our method on the DAVIS 2016 dataset [25], where the object categories may be different from the seen categories on PASCAL VOC. Based on this setting, we compare with the state-of-the-art methods for object segmentation via transfer learning, including approaches that use the NLP-based GloVe embeddings [24] and a decoupled network [10]. In addition, we show baseline results with and without the proposed iterative self-learning strategy to demonstrate its effectiveness. Second, we adopt the weakly-supervised setting on the YouTube-Objects dataset [28] and show that the proposed method performs favorably against the state-of-the-art algorithms in terms of visual quality and accuracy.

The contributions of this work are summarized as follows. First, we propose a framework for object segmentation in unlabeled videos through a self-supervised learning method. Second, we develop a joint formulation to mine useful segments while adapting the feature representations to the target videos. Third, we design a CNN module that can transfer knowledge from multiple seen categories in images to the unseen category in videos.

2 Related Work

Video Object Segmentation. Video object segmentation aims to separate foreground objects from the background. Conventional methods utilize object proposals [17,26,15] or graphical models [39,21], while recent approaches focus on learning CNN models from image sequences with frame-by-frame pixel-level

ground truth annotations to achieve the state-of-the-art performance [3,37,12]. For CNN-based methods, motion cues are usually used to effectively localize objects. Jain et al. [12] utilize a two-stream network by jointly considering appearance and motion information. The SegFlow method [3] further shows that jointly learning segmentation and optical flow in videos enhances both performance. Another line of research is to fine-tune the model based on the object mask in the first frame [1,14] and significantly improves the segmentation quality. However, in addition to annotations of the first frame in target videos [1,14], these methods require pre-training on videos with frame-by-frame pixel-level annotations [3,37] or bounding box ground truths [12] to obtain better foreground segmentation. In contrast, the proposed algorithm uses only a smaller number of existing annotations from the image dataset and transfers the feature representations to unlabeled videos for object segmentation. In addition, our method is flexible for the weakly-supervised learning setting, which cannot be achieved by the above approaches.

Object Segmentation in Weakly-supervised Videos. To reduce the need of pixel-level annotations, weakly-supervised methods have been developed to facilitate the segmentation process, where only class-level labels are required in videos. Numerous approaches are proposed to collect useful semantic segments by training segment-based classifiers [34] or ranking supervoxels [45]. However, these methods rely on the quality of generated segment proposals and may produce inaccurate results when taking low-quality segments as the input. Zhang et al. [44] propose to utilize object detectors integrated with object proposals to refine segmentations in videos. Furthermore, Tsai et al. [40] develop a co-segmentation framework by linking object tracklets from all the videos and improve the result. Recently, the SPFTN method [42] utilizes a self-paced learning scheme to fine-tune segmentation results from object proposals. Different from the above algorithms that only target on a pre-defined set of categories, our approach further extends this setting to videos without any labels for unseen object categories.

Transfer Learning for Object Recognition. Using cross-domain data for unsupervised learning has been explored in domain adaptation [30,7,23,6]. While domain adaptation methods make the assumption that the same categories are shared across different domains, transfer learning approaches focus on transferring knowledge between categories. Numerous transfer learning methods have been developed for object classification [38] and detection [18,9]. Similar efforts have been made for object segmentation. Hong et al. [10] propose a weakly-supervised semantic segmentation method by exploiting pixel-level annotations from different categories. Recently, Hu et al. [11] design a weighted transform function to transfer knowledge between the detected bounding boxes and instance segments. In this work, we share the similar motivation with [10] but remove the assumption of weak supervisions. To the best of our knowledge, this work is the first attempt for video object segmentation by transferring knowledge from annotated images to unlabeled videos between unshared categories.

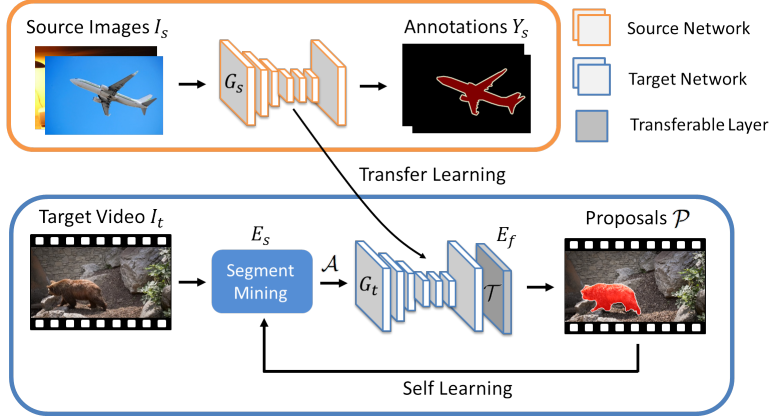


Fig. 1. Overview of the proposed algorithm. Given a set of source images \mathcal{I}_s with semantic segmentation annotations Y_s , we first train a source CNN model G_s . To predict object segmentations on a target video \mathcal{I}_t without knowing any annotations, we initialize the target network G_t from the parameters in G_s and perform adaptation via a transferable layer \mathcal{T} . We optimize the function E_s for selecting object-like segments \mathcal{A} from proposals \mathcal{P} and adapt feature representations in the CNN model via optimizing E_f . The entire self-learning process is performed via iteratively updating two energy functions to obtain the final segmentation results.

3 Algorithmic Overview

3.1 Overview of the Proposed Framework

We first describe the problem context of this work. Given a number of source images $\mathcal{I}_s = \{I_s^1, \dots, I_s^N\}$ with pixel-level semantic segmentation annotations $Y_s = \{y_s^1, \dots, y_s^N\}$ and the target sequence $\mathcal{I}_t = \{I_t^1, \dots, I_t^M\}$ without any labels, our objective is to develop a self-supervised learning algorithm that segments the object in \mathcal{I}_t by transferring knowledge from \mathcal{I}_s to \mathcal{I}_t . In this work, the object category in \mathcal{I}_t is allowed to be arbitrary. It can be either covered by or different from those in \mathcal{I}_s .

To this end, we propose a method with two components: 1) a ranking module for mining segment proposals, and 2) a CNN model for learning transferable feature representations. Fig. 1 illustrates these two components in the proposed framework. We first train a source CNN model G_s using \mathcal{I}_s and Y_s as the input and the desired output, respectively. Then we initialize the target network G_t from the parameters in G_s , where this target network can generate segment proposals \mathcal{P} on the target video \mathcal{I}_t . To find a set of object-like proposals among \mathcal{P} , we then develop an energy function to re-rank these proposals based on their objectness scores and mutual relationships. With the selected proposals that have higher object-like confidence, we further refine the feature representations in the target network. Since \mathcal{I}_s and \mathcal{I}_t may not share common object categories, we design a layer \mathcal{T} that enables cross-category knowledge transfer, and append it to

the target network. The entire process can be formulated as a joint optimization problem with the objective function as described below.

3.2 Objective Function

Our goal is to find high-quality segment proposals \mathcal{P} from the target video \mathcal{I}_t that can guide the network to learn feature representations \mathcal{F} for better segmenting the given video \mathcal{I}_t . We carry out this task by jointly optimizing an energy function E that accounts for segment proposals \mathcal{P} and features \mathcal{F} :

$$\max_{\mathcal{A}, \theta} E(\mathcal{I}_t, \mathcal{P}, \mathcal{F}; \mathcal{A}, \theta) = \max_{\mathcal{A}, \theta} E_s(\mathcal{P}, \mathcal{F}; \mathcal{A}) + E_f(\mathcal{I}_t, \mathcal{A}; \theta), \quad (1)$$

where E_s is the energy for selecting a set of high-quality segments \mathcal{A} from the proposals \mathcal{P} based on the features \mathcal{F} , while θ is the parameters of the CNN model that aims to optimize E_f and learn feature representations \mathcal{F} from the selected proposals \mathcal{A} . Details of each energy function and the optimization process are described in the following section.

4 Transferring Visual Information for Segmentation

In this section, we describe the details of the proposed energy functions for mining segments and learning feature representations, respectively. The segment mining step is formulated as a submodular optimization problem, while the feature learning process is completed through a CNN with a transferable module. After introducing both energy functions, we present an iterative optimization scheme to jointly maximize the objective (1).

4.1 Mining Segment Proposals

Given a target video \mathcal{I}_t , we can generate frame-by-frame object segmentations by applying the CNN model pre-trained on the source images \mathcal{I}_s . However, these segments may contain many false positives that do not well cover objects. Thus, we aim to select high-quality segments and eliminate noisy ones from the generated object segmentations. The challenging part lies in that there are no ground truth annotations in the target video, and thus we cannot train a classifier to guide the selection process.

Motivated by the co-segmentation method [40], we observe that high-quality segments have higher mutual relationships. As a result, we gather all the predicted segments from the target video and construct a graph to link each segment. We then formulate segment mining as a submodular optimization problem, aiming to select a subset of object-like segments that share higher similarities.

Graph Construction on Segments. We first feed the target video \mathcal{I}_t into the CNN model frame-by-frame and obtain a set of segment proposals \mathcal{P} , where each proposal is a connected-component in the predicted segmentation. Then

we construct a graph $G = (\mathcal{V}, \mathcal{E})$ on the set \mathcal{P} , where each vertex $v \in \mathcal{V}$ is a segment, and each edge $e \in \mathcal{E}$ models the pairwise relationship between two segments. Our goal is to find a subset \mathcal{A} within \mathcal{P} that contains proposals with higher object-like confidence.

Submodular Function. We design a submodular function to find segments that meet the following criteria: 1) objects from the same category share similar features, 2) a true object has a higher response from the output of the CNN model, and 3) an object usually moves differently from the background area in the video. Therefore, we formulate the objective function for selecting object-like segments by a facility location term \mathcal{H} [16] and a unary term \mathcal{U} . The former computes the similarity between the selected segments, while the latter estimates the probability of each selected segment being a true object. Both terms are defined based on the segment proposals \mathcal{P} and the adopted feature representation \mathcal{F} . First, we define the facility location term as:

$$\mathcal{H}(\mathcal{P}, \mathcal{F}; \mathcal{A}) = \sum_{i \in \mathcal{A}} \sum_{j \in \mathcal{V}} W(v_i, v_j) - \sum_{i \in \mathcal{A}} \phi_i, \quad (2)$$

where W denotes the pairwise relationship between a potential facility v_i and a vertex v_j , while ϕ_i is the cost to open a facility, which is fixed to a constant α . We define W as the similarity between two segments in order to encourage the submodular function to choose a facility v_i that is similar to v_j . To estimate this similarity, we represent each segment as a feature vector and compute their inner product of the two vectors. To form the feature vector for each segment, we draw feature maps from the CNN model (**conv1** to **conv5**) and perform the global average pooling on each segment. It is the adopted feature representation \mathcal{F} in this work. In addition to the facility location term, we employ a unary term to evaluate the quality of segments:

$$\mathcal{U}(\mathcal{P}, \mathcal{F}; \mathcal{A}) = \lambda_o \sum_{i \in \mathcal{A}} \Phi_o(i) + \lambda_m \sum_{i \in \mathcal{A}} \Phi_m(i), \quad (3)$$

where $\Phi_o(i)$ is the objectness score that measures the probability of segment i being a true object, and $\Phi_m(i)$ is the motion score that estimates the motion difference between segment i and the background region. λ_o and λ_m are the weights for the two terms, respectively.

The objectness score $\Phi_o(i)$ is calculated by averaging the probability map of the CNN output layer on all the pixels within the segment. For the motion score $\Phi_m(i)$, we first compute the optical flow [19] for two consecutive frames, and then we utilize the minimum barrier distance [33, 43] to convert the optical flow into a saliency map, where the larger distance represents a larger motion difference with respect to the background region.

Formulation for Segment Mining. Our goal is to find a subset \mathcal{A} within \mathcal{P} containing segments that are similar to each other and have higher object-like confidence. Therefore, we combine the facility location term \mathcal{H} and the unary

term \mathcal{U} as the energy E_s in (1):

$$E_s(\mathcal{P}, \mathcal{F}; \mathcal{A}) = \mathcal{H}(\mathcal{P}, \mathcal{F}; \mathcal{A}) + \mathcal{U}(\mathcal{P}, \mathcal{F}; \mathcal{A}). \quad (4)$$

We also note that the linear combination of two non-negative terms preserves the submodularity [46].

4.2 Learning Transferable Feature Representations

Given the selected set of object-like segment proposals, the ensuing task is to learn better feature representations based on these segments. To this end, we propose to use a CNN model fine-tuned on these segments via a self-learning scheme. However, since our target video may have a different set of object categories from those in the source domain, we further develop a transfer learning method where a transferable layer is augmented to the CNN model. With the proposed layer, our network is able to transfer knowledge from seen categories to the unseen category, without the need of any supervision in the target video.

Inspired by the observation that an unseen object category can be represented by a series of seen objects [29], we develop a transferable layer that approximates an unseen category as a linear combination of seen ones in terms of the output feature maps. In the following, we first present our CNN objective for learning the feature representations based on the selected segment proposals. Then we introduce the details of the proposed layer for transferring knowledge from the source domain to the target one.

Objective Function. Given the target video \mathcal{I}_t and selected segment proposals \mathcal{A} as described in Section 4.1, we use \mathcal{A} as our pseudo ground truths and optimize the target network G_t with parameters θ_g to obtain better feature representations that match the target video. Specifically, we define the energy function E_f in (1) as:

$$E_f(\mathcal{I}_t, \mathcal{A}; \theta_g, \theta_{\mathcal{T}}) = -\mathcal{L}(\mathcal{T}(G_t(\mathcal{I}_t)), \mathcal{A}), \quad (5)$$

where $\theta_{\mathcal{T}}$ is the parameters of the transferable layer \mathcal{T} and \mathcal{L} is the cross-entropy function to measure the loss between the network prediction $\mathcal{T}(G_t(\mathcal{I}_t))$ and the pseudo ground truth \mathcal{A} . We also note that, we use the minus sign for the loss function \mathcal{L} to match the maximization formulation in (1).

Learning Transferable Knowledge. Suppose there are C_s categories in the source domain, we aim to transfer a source network G_s pre-trained on the source images \mathcal{I}_s to the target video. To achieve this, we first initialize the target network G_t using the parameters in G_s . Given the target video \mathcal{I}_t , we can generate frame-wise feature maps $R = G_t(\mathcal{I}_t) = \{r_c\}_{c=1}^{C_s}$ through the network with C_s channels, where r_c is the output map of source category c . Since the target category is unknown, we then approximate the desired output map, r , for the unseen category as a linear combination of these seen categories through the proposed transferable layer \mathcal{T} :

$$r = \mathcal{T}(R) = \sum_{c=1}^{C_s} w_c r_c, \quad (6)$$

where w_c is the weight of the seen category c . Specifically, the proposed transferable layer \mathcal{T} can be performed via a 1×1 convolutional layer with C_s channels, in which the parameter of channel c in $\theta_{\mathcal{T}}$ corresponds to w_c .

Since w_c is not supervised by any annotations from the target video, the initialization of w_c is critical for obtaining a better combination of feature maps from the seen categories. Thus, we initialize w_c by calculating the similarity between each source category c and the target video. For each image in the source and target domains, we extract its feature maps from the **fc7** layer of the network and compute a 4096-dimensional feature vector on the predicted segment via global average pooling. By representing each image as a feature vector, we measure the similarity score between source and target images by their inner product. Finally, the initialized weight w_c^{init} for the category c can be obtained by averaging largest scores on each target frame with respect to the source images:

$$w_c^{init} = \frac{1}{|\mathcal{I}_t|} \sum_{i=1}^{|\mathcal{I}_t|} \max_j \langle \mathcal{F}_t^i, \mathcal{F}_{s,c}^j \rangle, \quad (7)$$

where $|\mathcal{I}_t|$ is the number of frames in the target video, $\mathcal{F}_t^i \in \mathbb{R}^{4096}$ is the feature vector of the i th frame of \mathcal{I}_t , and $\mathcal{F}_{s,c}^j \in \mathbb{R}^{4096}$ is the feature vector of the j th image of source category c .

4.3 Joint Formulation and Model Training

Based on the formulations to mine segments (4) and learn feature representations (5), we jointly solve the two objectives, i.e., E_s and E_f , in (1) by:

$$\begin{aligned} \max_{\mathcal{A}, \theta} E(\mathcal{I}_t, \mathcal{P}, \mathcal{F}; \mathcal{A}, \theta) &= \max_{\mathcal{A}, \theta} E_s(\mathcal{P}, \mathcal{F}; \mathcal{A}) + E_f(\mathcal{I}_t, \mathcal{A}; \theta) \\ &= \max_{\mathcal{A}, \theta_g, \theta_{\mathcal{T}}} [\mathcal{H}(\mathcal{P}, \mathcal{F}; \mathcal{A}) + \mathcal{U}(\mathcal{P}, \mathcal{F}; \mathcal{A})] - \mathcal{L}(\mathcal{T}(G_t(\mathcal{I}_t)), \mathcal{A}). \end{aligned} \quad (8)$$

To optimize (8), we decompose the process into two sub-problems: 1) solving a submodular function for segment mining to generate \mathcal{A} , and 2) training a CNN model that optimizes θ_g and $\theta_{\mathcal{T}}$ for learning transferable representations. We adopt an iterative procedure to alternately optimize the two sub-problems. The initialization strategy and the optimization of the two sub-problems are described below.

Initialization. We first pre-train a source network G_s on the PASCAL VOC training set [4] with 20 object categories. We then initialize the target network G_t from parameters in G_s and the transferable layer \mathcal{T} as described in Section 4.2. To obtain an initial set of segment proposals, we forward the target video \mathcal{I}_t to the target model G_t with \mathcal{T} and generate segments \mathcal{P} with their features \mathcal{F} .

Fix E_f and Optimize E_s . We first fix the network parameters θ and optimize \mathcal{A} in E_s of (8). We adopt a greedy algorithm similar to [40]. Starting from an

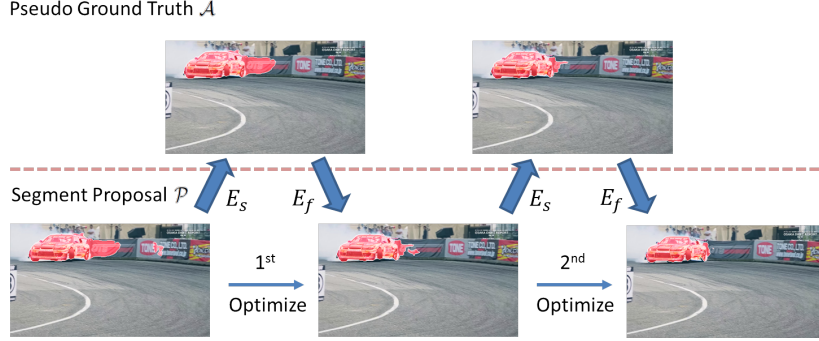


Fig. 2. Sample results of iteratively optimizing E_s and E_f . Starting from an initial set of proposals \mathcal{P} , we solve E_s to obtain object-like segments \mathcal{A} as our pseudo ground truths to optimize E_f . By iteratively updating both energy functions, our algorithm gradually improves the quality of \mathcal{P} and \mathcal{A} to obtain the final segmentation results.

empty set of \mathcal{A} , we add an initial element $a \in \mathcal{V} \setminus \mathcal{A}$ to \mathcal{A} that gives the largest energy gain. The process is then repeated and stops when one of the following conditions is satisfied: 1) the number of selected proposals reaches a threshold, i.e., $|\mathcal{A}| > N_{\mathcal{A}}$, and 2) the ratio of the energy gain between two rounds is below a threshold, i.e., $\mathcal{D}(\mathcal{A}^i) < \beta \cdot \mathcal{D}(\mathcal{A}^{i-1})$, where $\mathcal{D}(\mathcal{A}^i)$ stands for the energy gain, i.e., difference of E_s between two rounds during the optimization process, and β is the ratio.

Fix E_s and Optimize E_f . Once obtaining \mathcal{A} as the pseudo ground truths, we fix \mathcal{A} and optimize the network with the transferable layer, i.e., θ_g and $\theta_{\mathcal{T}}$, in E_f of (8). We alter the problem to a task that minimizes the network loss \mathcal{L} in an end-to-end fashion, jointly for θ_g and $\theta_{\mathcal{T}}$ using the SGD method.

Iterative Optimization. To obtain the final \mathcal{A} , θ_g and $\theta_{\mathcal{T}}$, instead of directly solving (8) for optimal solutions, we solve it via an iterative updating scheme between E_s and E_f until convergence. In practice, we measure the intersection-over-union (IoU) of selected segment proposals between two iterations. The optimization process ends when the IoU is larger than a threshold (e.g., 90%), showing that the set of \mathcal{A} becomes stable. Fig. 2 shows one example of gradually improved \mathcal{P} and \mathcal{A} via iteratively updating E_s and E_f . The overall process is summarized in Algorithm 1.

5 Experimental Results

In this section, we first present implementation details of the proposed method, and then we show experimental results on numerous benchmark datasets. In addition, ablation studies for various components in the algorithm are conducted. The source code and trained models will be made available to the public. More results are presented in the supplementary material.

Algorithm 1 Unseen Object Segmentation

Source Image: \mathcal{I}_s, Y_s
Target Video: \mathcal{I}_t
Initialization: pre-trained G_s on source inputs, $G_t \leftarrow G_s, w_c^{init}$ via (7)
 $(\mathcal{P}, \mathcal{F}) \leftarrow \mathcal{T}(G_t(\mathcal{I}_t))$
while \mathcal{P} not converged **do**
 $\mathcal{A}^0 \leftarrow \emptyset, i \leftarrow 1$
loop
 $a^* = \arg \max_{\{\mathcal{A}^i \in \mathcal{V}\}} E_s(\mathcal{P}, \mathcal{F}; \mathcal{A}^i)$, where $\mathcal{A}^i \leftarrow \mathcal{A}^{i-1} \cup a, a \in \mathcal{V} \setminus \mathcal{A}$
if $|\mathcal{A}| > N_{\mathcal{A}}$ or $\mathcal{D}(\mathcal{A}^i) < \beta \cdot \mathcal{D}(\mathcal{A}^{i-1})$ when $i \geq 2$ **then**
break
end if
 $\mathcal{A}^i \leftarrow \mathcal{A}^{i-1} \cup a^*, i \leftarrow i + 1$
end loop
 $\mathcal{A} \leftarrow \mathcal{A}^i$
Optimize $E_f: (\theta_g, \theta_T) \leftarrow \min \mathcal{L}(\mathcal{T}(G_t(\mathcal{I}_t)), \mathcal{A})$
 $(\mathcal{P}, \mathcal{F}) \leftarrow \mathcal{T}(G_t(\mathcal{I}_t))$
end while
Output: object segmentation \mathcal{P} of \mathcal{I}_t

5.1 Implementation Details

In the submodular function for segment mining, we set $\alpha = 1$ for the facility location term in (2), and $\lambda_o = 20$, $\lambda_m = 35$ for the unary term in (3). During the submodular optimization in (4), we use $N_{\mathcal{A}} = 0.8 \cdot |\mathcal{P}|$ and $\beta = 0.8$. All the parameters are fixed in all the experiments. For training the CNN model in (5), we employ various fully convolutional networks (FCNs) [20] including the VGG-16 [32] and ResNet-101 [8] architectures for both the source and target networks using the Caffe library. The learning rate, momentum and batch size are set as 10^{-14} , 0.99, and 1, respectively. To further refine segmentation results, we average the responses from the CNN output and a motion prior that is already computed in the motion term of (3) to account for both the appearance and temporal information.

5.2 DAVIS Dataset

We first conduct experiments on the DAVIS 2016 benchmark dataset [25]. Since our goal is to transfer the knowledge from seen categories in images to unseen objects in the video, we manually select all the videos with object categories that are different from the 20 categories in the PASCAL VOC dataset. In the following, we first conduct ablation studies and experiments to validate the proposed method. Second, we show that our algorithm can be applied under various settings on the entire set of the DAVIS 2016 dataset.

Impact of the Motion Term. One critical component of our framework is to mine useful segments for the further CNN model training step. In the submodular function of (3), we incorporate a motion term that accounts for object

Table 1. IoU of the selected segments with different weights of the motion term on the DAVIS dataset.

λ_m	0	5	15	25	35	45
Avg. IoU	57.2	57.4	60.5	60.6	61.0	60.3

Table 2. Results on the DAVIS 2016 dataset with categories excluded from the PASCAL VOC dataset.

Methods	bear	bswan	camel	eleph	goat	malw	rhino	Avg.
CVOS [35]	86.4	42.2	85.0	49.4	7.4	24.5	52.0	49.6
MSG [27]	85.1	52.6	75.6	68.9	73.5	4.5	90.2	64.3
FST [22]	89.8	73.2	56.2	82.4	55.4	8.7	77.6	63.3
NLC [5]	90.7	87.5	76.8	51.8	1.0	76.1	68.2	64.6
LMP [36]	69.8	50.9	78.3	78.9	75.1	38.5	76.8	66.9
TransferNet [10]	73.7	83.4	65.5	76.1	78.1	17.9	42.4	62.4
Ours (GloVe)	82.6	67.2	68.8	61.2	70.4	64.7	32.0	63.8
Ours (init)	80.3	75.6	70.9	70.4	83.1	40.9	57.7	68.4
Ours (opt)	88.8	80.6	68.6	71.8	82.4	43.8	67.3	71.9
Ours (final)	89.8	76.7	72.0	73.8	83.3	41.6	71.0	72.6
ARP [15]	92	88.1	90.3	84.2	77.6	58.3	88.4	82.7
FSEG [12]	91.5	89.5	76.4	86.2	84.1	83.3	77.6	84.1
Ours (ResNet)	91.8	90.3	77.5	85.7	84.8	84.9	86.0	85.9

movements in the video. To validate its effectiveness, we fix the weight $\lambda_o = 20$ for the appearance and vary the weight λ_m for the motion term. In Table 1, we show the IoU of the selected segment proposals via solving (4) under various values of λ_m . The results show that the IoU is gradually improved when increasing the motion weight, which indicates that the quality of selected segments becomes better, and hence we use $\lambda_m = 35$ in all the following experiments.

Ablation Study. In the middle group of Table 2, we show the final segmentation results of our method using VGG-16 architecture with various baselines and settings. We first present a baseline method that uses the GloVe embeddings [24] to initialize weights, i.e., the similarity between two categories, of the transferable layer. Since the GloVe is not learned in the image domain between categories, the initialized weights may not reflect the true relationships between the seen and unseen categories, and hence the results are worse than the proposed method for initializing the transferable layer.

Furthermore, we show results at different stages, including using the model with initialization before optimizing (8), after optimization, and the final result with motion refinement. After the optimization, the IoU is improved in 5 out of 7 videos, which shows the effectiveness of the proposed self-learning scheme without using any annotations in the target video.

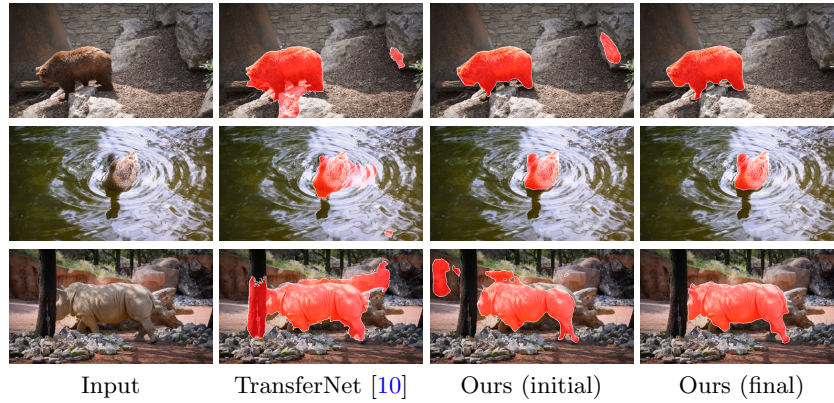


Fig. 3. Sample results on the DAVIS dataset with categories excluded from the PASCAL VOC dataset. We show that our final results are more accurate in details than the TransferNet [10] and with the noisy segments removed from the initial results.

Overall Comparisons. In Table 2, we show the comparisons between our method and the state-of-the-art approaches. We first demonstrate the performance of our method using VGG-16 architecture. The work closest in the scope to the proposed framework is the method [10] that transfers the knowledge between two image domains with mutually exclusive categories in a weakly-supervised setting. To compare with this approach, we use the authors’ public implementation and train the models with the same setting as our method. We show that our algorithm achieves better IoUs in 5 out of 7 videos and improves the overall IoU by 10.2% on average. We also note that our model with initialization already performs favorably against [10], which demonstrates that the proposed transferable layer is effective in learning knowledge from seen categories to unseen ones. Visual comparisons are presented in Fig. 3.

In addition, we present more results of video object segmentation methods in Table 2 and show that the proposed algorithm achieves better performance. Different from existing approaches that rely on long-term trajectory [35,27] or motion saliency [22,5] to localize foreground objects, we use the proposed self-learning framework to segment unseen object categories via transfer learning. We note that the proposed method performs better than the CNN-based model [36] that utilizes synthetic videos with pixel-wise segmentation annotations.

We further employ the stronger ResNet-101 architecture and compare with state-of-the-art unsupervised video object segmentation methods. In the bottom group of Table 2, we show that our approach performs better than FSEG [12] using the same architecture and training data from PASCAL VOC, i.e., the setting of the appearance stream in FSEG [12]. In addition, compared to ARP [15] that adopts a non-learning based framework via proposal post-processing and is specifically designed for video object segmentation, our algorithm performs better and is flexible under various settings such as using weakly-supervised signals.

Results on the Entire DAVIS 2016 Dataset. In addition to performing object segmentation on unseen object categories, our method can adapt to the weakly-supervised setting by simply initializing the weights in the transferable layer as a one-hot vector, where only the known category is set to 1 and the others are 0. We evaluate this setting on the DAVIS 2016 dataset with categories shared in the PASCAL VOC dataset. Note that, we still adopt the unsupervised setting for the unseen categories. The results on the entire DAVIS 2016 dataset are shown in Table 3. In comparison with a recent weakly-supervised method [42] and the baseline model [20] (our initial result), our approach addresses the transfer learning problem and outperforms their methods by 6.5% and 6.1%, respectively.

Although the same categories are shared between the source and target domains in this setting, we can still assume that the object category is unknown in the target video. Under this fully unsupervised setting without using any pixel-wise annotations in videos during training, we show that our method improves the results of FSEG [12] and other unsupervised algorithms [22,5]. Sample results are presented in Fig. 4.

Table 3. Results on the entire DAVIS 2016 dataset.

Methods	Weak Supervision			No Supervision			
	SPFTN [42]	FCN [20]	Ours	FST [22]	NLC [5]	FSEG [12]	Ours
Avg. IoU	61.2	61.6	67.7	57.5	64.1	64.7	66.5

5.3 YouTube-Objects Dataset

We evaluate our method on the YouTube-Objects dataset [28] with annotations provided by [13] for 126 videos. Since this dataset contains 10 object categories that are shared with the PASCAL VOC dataset, we conduct experiments using the weakly-supervised setting. In Table 4, we compare our method with the state-of-the-art algorithms that use the class-level weak supervision. With the VGG-16 architecture, the proposed framework performs well in 6 out of 10 categories and achieves the best IoU on average. Compared to the baseline FCN model [20] used in our algorithm, there is a performance gain of 9%. In addition, while existing methods rely on training the segment classifier [34], integrating object proposals with detectors [44], co-segmentation via modeling relationships between videos [40], or self-paced fine-tuning [42], the proposed method utilizes a self-learning scheme to achieve better segmentation results. With the ResNet-101 architecture, we compare our method with DeepLab [2] and FSEG [12]. We show that the proposed method improves the performance in 6 out of 10 categories and achieves the best averaged IoU.

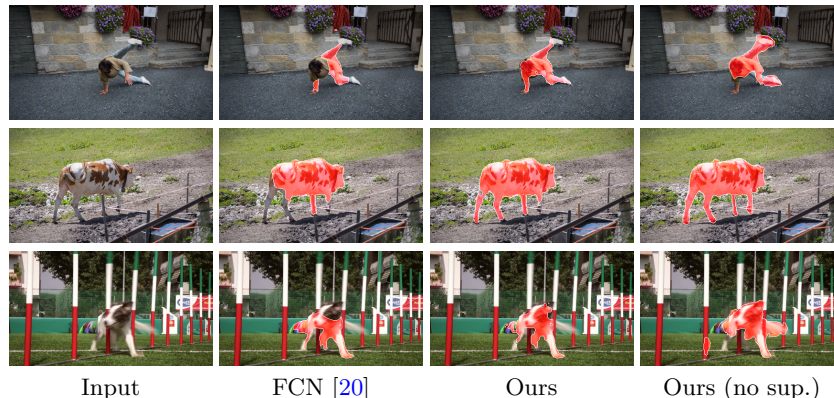


Fig. 4. Segmentation results on the DAVIS dataset with categories shared in the PASCAL VOC dataset. We show that both of our results with and without supervision have more complete object segmentations than the baseline FCN model [20] (our initial result) that uses the weak supervision.

Table 4. Results on the YouTube-Objects dataset.

Methods	aero	bird	boat	car	cat	cow	dog	horse	mbike	train	Avg.
DSA [34]	17.8	19.8	22.5	38.3	23.6	26.8	23.7	14.0	12.5	40.4	23.9
FCN [20]	68.3	65.7	55.7	76.6	52.3	50.4	55.6	52.6	35.7	55.9	56.9
DET [44]	72.4	66.6	43.0	58.9	36.4	58.2	48.7	49.6	41.4	49.3	52.4
CoSeg [40]	69.3	76.1	57.2	70.4	67.7	59.7	64.2	57.1	44.1	57.9	62.3
SPFTN [42]	81.1	68.8	63.4	73.8	59.7	64.5	63.4	58.2	52.4	45.5	63.1
Ours (VGG)	74.6	65.3	66.9	79.5	64.2	68.3	67.3	61.7	51.5	59.4	65.9
DeepLab [2]	80.6	67.8	66.9	73.3	55.3	61.8	63.9	45.5	54.7	56.4	62.6
FSEG [12]	83.4	60.9	72.6	74.5	68.0	69.6	69.1	62.8	61.9	62.8	68.6
Ours (ResNet)	83.5	76.4	70.0	75.3	65.9	69.7	71.6	54.7	63.8	58.7	69.0

6 Concluding Remarks

In this paper, we propose a self-learning framework to segment objects in unlabeled videos. By utilizing existing annotations in images, we design a model to adapt seen object categories from source images to the target video. The entire process is decomposed into two sub-problems: 1) a segment mining module to select object-like proposals, and 2) a CNN model with a transferable layer that adapts feature representations for target videos. To optimize the proposed formulation, we adopt an iterative scheme to obtain final solutions. Extensive experiments and ablation study show the effectiveness of the proposed algorithm against other state-of-the-art methods on numerous benchmark datasets.

Acknowledgments. This work is supported in part by Ministry of Science and Technology under grants MOST 105-2221-E-001-030-MY2 and MOST 107-2628-E-001-005-MY3.

References

1. Caelles, S., Maninis, K.K., Pont-Tuset, J., Leal-Taixé, L., Cremers, D., Gool, L.V.: One-shot video object segmentation. In: CVPR (2017) [3](#)
2. 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) [13](#), [14](#)
3. Cheng, J., Tsai, Y.H., Wang, S., Yang, M.H.: Segflow: Joint learning for video object segmentation and optical flow. In: ICCV (2017) [3](#)
4. Everingham, M., Gool, L.J.V., Williams, C.K.I., Winn, J.M., Zisserman, A.: The pascal visual object classes (voc) challenge. IJCV **88**(2), 303–338 (2010) [2](#), [8](#)
5. Faktor, A., Irani, M.: Video segmentation by non-local consensus voting. In: BMVC (2014) [11](#), [12](#), [13](#)
6. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. In: ICML (2015) [3](#)
7. Gopalan, R., Li, R., Chellappa, R.: Domain adaptation for object recognition: An unsupervised approach. In: ICCV (2011) [3](#)
8. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) [10](#)
9. Hoffman, J., Guadarrama, S., Tzeng, E.S., Hu, R., Donahue, J., Girshick, R., Darrell, T., Saenko, K.: Lsda: Large scale detection through adaptation. In: NIPS (2014) [3](#)
10. Hong, S., Oh, J., Lee, H., Han, B.: Learning transferrable knowledge for semantic segmentation with deep convolutional neural network. In: CVPR (2016) [2](#), [3](#), [11](#), [12](#)
11. Hu, R., Dollár, P., He, K., Darrell, T., Girshick, R.: Learning to segment every thing. arXiv:1711.10370 (2017) [3](#)
12. Jain, S., Xiong, B., Grauman, K.: Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. In: CVPR (2017) [3](#), [11](#), [12](#), [13](#), [14](#)
13. Jain, S.D., Grauman, K.: Supervoxel-consistent foreground propagation in video. In: ECCV (2014) [13](#)
14. Khoreva, A., Perazzi, F., Benenson, R., Schiele, B., Sorkine-Hornung, A.: Learning video object segmentation from static images. In: CVPR (2017) [3](#)
15. Koh, Y.J., Kim, C.S.: Primary object segmentation in videos based on region augmentation and reduction. In: CVPR (2017) [2](#), [11](#), [12](#)
16. Lazic, N., Givoni, I., Frey, B., Aarabi, P.: Floss: Facility location for subspace segmentation. In: ICCV (2009) [6](#)
17. Lee, Y.J., Kim, J., Grauman, K.: Key-segments for video object segmentation. In: ICCV (2011) [2](#)
18. Lim, J.J., Salakhutdinov, R., Torralba, A.: Transfer learning by borrowing examples for multiclass object detection. In: NIPS (2011) [3](#)
19. Liu, C.: Beyond pixels: Exploring new representations and applications for motion analysis. PhD thesis, MIT (2009) [6](#)
20. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: CVPR (2015) [10](#), [13](#), [14](#)
21. Márki, N., Perazzi, F., Wang, O., Sorkine-Hornung, A.: Bilateral space video segmentation. In: CVPR (2016) [2](#)
22. Papazoglou, A., Ferrari, V.: Fast object segmentation in unconstrained video. In: ICCV (2013) [11](#), [12](#), [13](#)

23. Patricia, N., Caputo, B.: Learning to learn, from transfer learning to domain adaptation: A unifying perspective. In: CVPR (2014) [3](#)
24. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: EMNLP. pp. 1532–1543 (2014) [2](#), [11](#)
25. 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) [2](#), [10](#)
26. Perazzi, F., Wang, O., Gross, M., Sorkine-Hornung, A.: Fully connected object proposals for video segmentation. In: CVPR (2015) [2](#)
27. P.Ochs, T.Brox: Object segmentation in video: a hierarchical variational approach for turning point trajectories into dense regions. In: ICCV (2011) [11](#), [12](#)
28. Prest, A., Leistner, C., Civera, J., Schmid, C., Ferrari, V.: Learning object class detectors from weakly annotated video. In: CVPR (2012) [2](#), [13](#)
29. Rochan, M., Wang, Y.: Weakly supervised localization of novel objects using appearance transfer. In: CVPR (2015) [7](#)
30. Saenko, K., Kulis, B., Fritz, M., Darrell, T.: Adapting visual category models to new domains. In: ECCV (2010) [3](#)
31. Saleh, F.S., Aliakbarian, M.S., Salzmann, M., Petersson, L., Alvarez, J.M.: Bringing background into the foreground: Making all classes equal in weakly-supervised video semantic segmentation. In: ICCV (2017) [1](#)
32. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. CoRR **abs/1409.1556**, 1187–1200 (2014) [10](#)
33. Strand, R., Ciesielski, K.C., Malmberg, F., Saha, P.K.: The minimum barrier distance. CVIU **117**(4), 429–437 (2013) [6](#)
34. Tang, K., Sukthankar, R., Yagnik, J., Fei-Fei, L.: Discriminative segment annotation in weakly labeled video. In: CVPR (2013) [3](#), [13](#), [14](#)
35. Taylor, B., Karasev, V., Soatto, S.: Causal video object segmentation from persistence of occlusions. In: CVPR (2015) [11](#), [12](#)
36. Tokmakov, P., Alahari, K., Schmid, C.: Learning motion patterns in videos. In: CVPR (2017) [11](#), [12](#)
37. Tokmakov, P., Alahari, K., Schmid, C.: Learning video object segmentation with visual memory. In: ICCV (2017) [3](#)
38. Tommasi, T., Orabona, F., Caputo, B.: Learning categories from few examples with multi model knowledge transfer. PAMI **36**, 928941 (2014) [3](#)
39. Tsai, Y.H., Yang, M.H., Black, M.J.: Video segmentation via object flow. In: CVPR (2016) [2](#)
40. Tsai, Y.H., Zhong, G., Yang, M.H.: Semantic co-segmentation in videos. In: ECCV (2016) [1](#), [3](#), [5](#), [8](#), [13](#), [14](#)
41. Yan, Y., Xu, C., Cai, D., Corso, J.J.: Weakly supervised actor-action segmentation via robust multi-task ranking. In: CVPR (2017) [1](#)
42. Zhang, D., Yang, L., Meng, D., Xu, D., Han, J.: Spftn: A self-paced fine-tuning network for segmenting objects in weakly labelled videos. In: CVPR (2017) [1](#), [3](#), [13](#), [14](#)
43. Zhang, J., Sclaroff, S., Lin, Z., Shen, X., Price, B., Mech, R.: Minimum barrier salient object detection at 80 fps. In: ICCV (2015) [6](#)
44. Zhang, Y., Chen, X., Li, J., Wang, C., Xia, C.: Semantic object segmentation via detection in weakly labeled video. In: CVPR (2015) [3](#), [13](#), [14](#)
45. Zhong, G., Tsai, Y.H., Yang, M.H.: Weakly-supervised video scene co-parsing. In: ACCV (2016) [3](#)
46. Zhu, F., Jiang, Z., Shao, L.: Submodular object recognition. In: CVPR (2014) [7](#)