

Weakly Supervised Instance Segmentation for Videos with Temporal Mask Consistency

Qing Liu*
Johns Hopkins University
qingliu@jhu.edu

Vignesh Ramanathan
Facebook
vigneshr@fb.com

Dhruv Mahajan
Facebook
dhruvm@fb.com

Alan Yuille
Johns Hopkins University
alan.l.yuille@gmail.com

Zhenheng Yang
Facebook
zhenheny@gmail.com

Abstract

Weakly supervised instance segmentation reduces the cost of annotations required to train models. However, existing approaches which rely only on image-level class labels predominantly suffer from errors due to (a) partial segmentation of objects and (b) missing object predictions. We show that these issues can be better addressed by training with weakly labeled videos instead of images. In videos, motion and temporal consistency of predictions across frames provide complementary signals which can help segmentation. We are the first to explore the use of these video signals to tackle weakly supervised instance segmentation. We propose two ways to leverage this information in our model. First, we adapt inter-pixel relation network (IRN) [6] to effectively incorporate motion information during training. Second, we introduce a new MaskConsist module, which addresses the problem of missing object instances by transferring stable predictions between neighboring frames during training. We demonstrate that both approaches together improve the instance segmentation metric AP_{50} on video frames of two datasets: Youtube-VIS and Cityscapes by 5% and 3% respectively.

1. Introduction

Instance segmentation is a challenging task, where all object instances in an image have to be detected and segmented. This task has seen rapid progress in recent years [17, 34, 10], partly due to the availability of large datasets like COCO [32]. However, it can be forbiddingly expensive to build datasets at this scale for a new domain of images or videos, since segmentation boundaries have to be annotated for every object in an image.

Alternatively, weak labels like classification labels can be used to train instance segmentation models [61, 11, 62, 6,

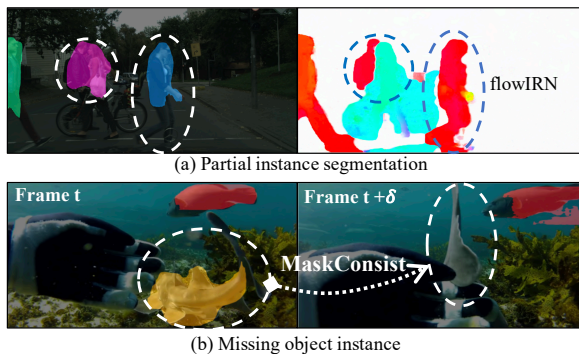


Figure 1. Two types of error for IRN [6] trained with still images: (a) partial segmentation and (b) missing instance. We observe optical flow is able to capture pixels of the same instance better (circles in (a)) and we propose flowIRN to model this information. In (b), a fish is missed on one frame. We propose MaskConsist to leverage temporal consistency and transfer stable mask predictions to neighboring frames during training.

27, 16, 44, 7]. While weak labels are significantly cheaper to annotate, training weakly supervised models can be far more challenging. They typically suffer from two sources of error: (a) *partial instance segmentation* and (b) *missing object instances*, as shown in Fig. 1. Weakly supervised methods often identify only the most discriminative object regions that help predict the class label. This results in *partial segmentation of objects*, as shown in Fig. 1(a). For instance, the recent work on weakly supervised instance segmentation IRN [6] relies on class activation maps (CAMs) [60], which suffer from this issue as also observed in other works [26, 55, 59]. Further, CAMs do not differentiate between overlapping instances of the same class. It can also *miss object instances* when multiple instances are present in an image, as shown in Fig. 1(b). In particular, an instance could be segmented in one image but not in another image where it is occluded or its pose alters.

Interestingly, these issues are less severe in videos, where object motion provides an additional signal for in-

*This work is done during Qing Liu’s internship at Facebook.

stance segmentation. As shown in Fig. 1, optical flow in a video is tightly coupled with instance segmentation masks. This is unsurprising since pixels belonging to the same (rigid) object move together and have similar flow vectors. We incorporate such video signals to train weakly supervised instance segmentation models, in contrast to existing methods [6, 27, 44, 7] only targeted at images.

Typical weakly supervised approaches involve two steps: (a) generating pseudo-labels, comprising noisy instance segmentation masks consistent with the weak class labels, and (b) training a supervised model like Mask R-CNN based on these pseudo-labels. We leverage video information in both stages. In the first step, we modify IRN to assign similar labels to pixels with similar motion. This helps in addressing the problem of partial segmentation. We refer to the modified IRN as *flowIRN*. In the second step, we introduce a new module called *MaskConsist*, which counters the problem of missing instances by leveraging temporal consistency between objects across consecutive frames. It matches prediction between neighboring frames and transfers the stable predictions to obtain additional pseudo-labels missed by *flowIRN* during training. This is a generic module that can be used in combination with any weakly supervised segmentation methods as we show in our experiments.

To the best of our knowledge, we are the first work to utilize temporal consistency between frames to train a weakly supervised instance segmentation model for videos. We show that this leads to more than 5% and 3% improvement in average precision compared to image-centric methods, like IRN, on video frames from two challenging video datasets: Youtube-VIS (YTVIS) [58] and Cityscapes [12], respectively. We also observe similar gains on the recently introduced video instance segmentation task [58] in YTVIS.

2. Related Work

Different types of weak supervision have been used in the past for semantic segmentation: bounding boxes [13, 39, 25, 46], scribbles [31, 51, 48], and image-level class labels [26, 24, 19, 56, 22, 28, 45, 47]. Similarly, for instance segmentation, image-level [61, 11, 62, 6, 27, 16, 44, 7] and bounding box supervision [25, 20] have been explored. In this work, we focus on only using class labels for weakly supervised instance segmentation.

Weakly supervised semantic segmentation: Most weakly supervised semantic segmentation approaches rely on class attention maps (CAMs) [60] to provide noisy pseudo-labels as supervision [26, 22, 45, 47]. Sun *et al.* [47] used co-attention maps generated from image pairs to train the semantic segmentation network. Another line of work leverages motion and temporal consistency in videos [50, 49, 43, 18, 29, 54] to learn more robust representation. For instance, frame-to-frame (F2F) [29] used optical flow to warp CAMs from neighboring frames and aggregated the warped CAMs to obtain more robust pseudo-labels.

Weakly supervised instance segmentation: For training instance segmentation models with bounding box supervision, Hsu *et al.* [20] proposed a bounding box tightness constraint and multiple instance learning (MIL) based objective. Another line of work that only uses class labels extracts semantic responses from CAMs or other attention maps and then combines them with object proposals [42, 41] to generate instance segmentation masks [61, 11, 62, 27, 44]. However, these methods’ performance heavily depends on the quality of proposals used, which are mostly pre-trained on other datasets. Shen *et al.* [44] extracted attention maps from a detection network and then jointly learn the detection and segmentation networks in a cyclic manner. Arun *et al.* [7] proposed a conditional network to model the noise in weak supervision and combined it with object proposals to generate instance masks. The first end-to-end network (IRN) [6] was proposed by Ahn *et al.* to directly predict instance offset and semantic boundary which were combined with CAMs to generate instance mask predictions. Our method adapts [6] for the first step of training and combines it with a novel *MaskConsist* module. However, other weakly supervised can also be integrated into our framework if code is available.

Segmentation in videos: A series of approaches have emerged for segmentation in videos [9, 21, 8, 38, 30]. Some works proposed to leverage the video consistency [52, 53, 36, 37]. Recently, Yang *et al.* [58] extended the traditional instance segmentation task from images to videos and proposed Video Instance Segmentation task (VIS). VIS aims to simultaneously segment and track all object instances in the video. Every pixel is labeled with a class label and an instance track-ID. MaskTrack [58] added a tracking-head to Mask R-CNN [17] to build a new model for this task. Bertasius *et al.* [9] improved MaskTrack by proposing a mask propagation head. This head propagated instance features across frames in a clip to get more stable predictions. To the best of our knowledge, there has been no work that has explored weakly supervised learning for the video instance segmentation task. We evaluate our method on this task by combining it with a simple tracking approach.

3. Approach

We first introduce preliminaries of inter-pixel relation network (IRN) [6] and extend it to incorporate video information, resulting in *flowIRN*. Next, we introduce *MaskConsist* which enforces temporal consistency in predictions across successive frames. Our framework has a 2-stage training process: (1) train *flowIRN* and (2) use masks generated by *flowIRN* on the training frames as supervision to train the *MaskConsist* model, as shown in Fig. 2.

3.1. Preliminaries of IRN

IRN [6] extracts inter-pixel relations from Class Attention Maps (CAMs) and uses it to infer instance locations

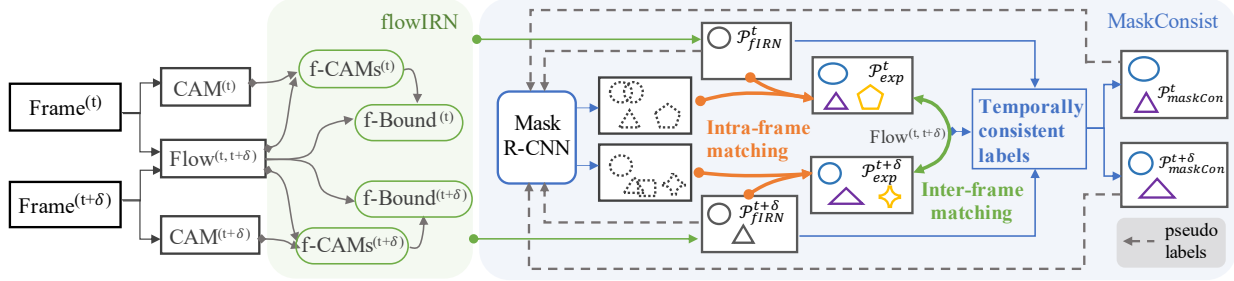


Figure 2. Our pipeline mainly consists of two modules: flowIRN and MaskConsist. FlowIRN adapts IRN [6] by incorporating optical flow to modify CAMs (f-CAMs), as well as introducing a new loss function: flow-boundary loss (f-Bound loss). MaskConsist matches the predictions from two successive frames and transfers high-quality predictions from one frame as pseudo-labels to another. It has three components: intra-frame matching, inter-frame matching and temporally consistent labels, shown in orange, green and blue, respectively. First, flowIRN is trained with frame-level class labels. Next, MaskConsist is trained with the pseudo-labels generated by flowIRN.

and class boundaries. For a given image, CAMs provide pixel-level scores for each class that are then converted to class labels. Every pixel is assigned the label corresponding to the highest class activation score at the pixel, if this score is above a foreground threshold. Otherwise, it is assigned the background label.

IRN is a network with two branches that predict (a) a per-pixel *displacement vector* pointing towards the center of the instance containing the pixel and (b) a per-pixel *boundary likelihood* indicating if a pixel lies on the boundary of an object or not. Since the model is weakly supervised, neither displacement nor boundary labels are available during training. Instead, IRN introduces losses that enforce constraints on displacement and boundary predictions based on the foreground/background labels inferred from CAMs.

During inference, a two-step procedure is used to obtain instance segmentation masks. First, all pixels with displacement vectors pointing towards the same centroid are grouped together to obtain per-pixel instance labels. However, these predictions tend to be noisy. In the second step, the predictions are refined using a pairwise affinity term α . For two pixels i and j ,

$$\alpha_{i,j} = 1 - \max_{k \in \Pi_{i,j}} \mathcal{B}(k), \quad (1)$$

where $\mathcal{B}(k)$ is the boundary likelihood predicted by IRN for pixel k , and $\Pi_{i,j}$ is the set of pixels lying on the line connecting i and j . If two pixels are separated by an object boundary, at least one pixel on the line connecting them should belong to this boundary. This results in low affinity between the two pixels. Conversely, the affinity would be high for pixels which are part of the same instance. In IRN, the affinity term is used to define the transition probability for a random walk algorithm that smooths the final per-pixel instance and class label assignments.

3.2. FlowIRN Module

We introduce flowIRN which improves IRN by incorporating optical flow information in two components, flow-amplified CAMs and flow-boundary loss.

Flow-Amplified CAMs: We observed that CAMs iden-

tify only the discriminative regions of an object (like the face of an animal) but often miss other regions corresponding to the object. This has been noted in previous works as well [26, 55, 59]. Since the objects of interest in a video are usually moving foreground objects, we address this issue by first amplifying CAMs in regions where large motion is observed. More specifically, given the estimated optical flow $\mathcal{F} \in \mathbb{R}^{H \times W \times 2}$ for the current frame, we replace CAMs used in IRN with:

$$\text{f-CAM}_c(x) = \text{CAM}_c(x) \times A^{\mathbb{I}(\|\mathcal{F}(x)\|_2 > T)}, \quad (2)$$

where A is an amplification coefficient and T is a flow magnitude threshold. This operation is applied to CAMs of all classes equally, preserving the relative ordering of class scores. Class labels obtained from CAMs are not flipped; only foreground and background assignments are affected.

Flow-boundary loss: In IRN, boundary prediction is supervised by the pseudo segmentation labels from CAMs, which does not distinguish instances of the same class, particularly overlapping instances. However, in videos, optical flow could disambiguate such instances, since pixels of the same rigid object move together and have consistent motion. Hence, we use spatial gradient of optical flow to identify if two pixels are from the same object instance or not. Points from the same object can be from different depths relative to the camera, and might not have the same optical flow. In practice, we observed that the gradient is more robust to this depth change. We explain this in detail in the appendix. We use the affinity term from Eq. 1 to define a new flow-boundary loss:

$$\mathcal{L}_{\mathcal{F}}^{\mathcal{B}} = \sum_{j \in \mathcal{N}_i} \|\mathcal{F}'(i) - \mathcal{F}'(j)\| \alpha_{i,j} + \lambda |1 - \alpha_{i,j}|, \quad (3)$$

where $\mathcal{F}'(i)$ is a two-dimensional vector denoting the gradient of the optical flow at a pixel i with respect to its spatial co-ordinates (x_i, y_i) respective. \mathcal{N}_i is a small pixel neighborhood around i as defined in IRN, and λ is the regularization parameter. The first term implies that pixels with similar flow-gradients could have high-affinity (belonging to the same instance), while pixels with different flow-gradients

should have low-affinity (belonging to different instances). The second term is used to regularize the loss and prevent the trivial solution of α being 0 constantly. We train flowIRN with the above loss and the original losses in IRN.

3.3. MaskConsist

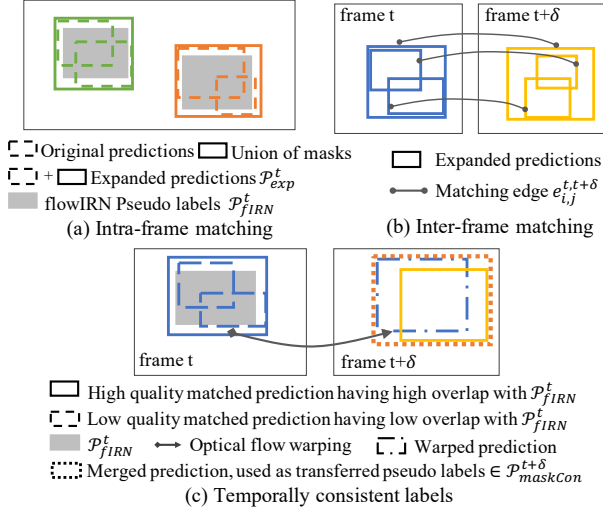


Figure 3. Three steps of MaskConsist module. (a) Intra-frame matching expands the original predictions of current Mask R-CNN by merging highly overlapping predictions. (b) Inter-frame matching identifies one-to-one matching between predictions across frames. (c) Temporally consistent labels transfer matched predictions from one frame to another after warping with optical flow.

The instance-level masks generated by flowIRN can now be used as pseudo-labels to train a fully supervised instance segmentation model, like Mask R-CNN. In practice, this yields better performance on the validation set than the original flowIRN model. However, the pseudo-labels generated by flowIRN can miss object instances on some frames if both CAMs and optical-flow couldn’t identify them.

MaskConsist solves this by transferring “high-quality” mask predictions from a neighboring frame as new pseudo-labels to the current frame while training a Mask R-CNN. At each training iteration, we train the network with a pair of neighboring frames t and $t + \delta$ in the video. In addition to the pseudo-labels from flowIRN, a subset of predictions by the current Mask R-CNN on $t + \delta$ are used as additional pseudo-labels for t and vice-versa. Predictions are transferred only if (a) they are temporally stable and (b) overlap with existing pseudo-labels from flowIRN. This avoids false-negatives by MaskConsist. MaskConsist contains three steps: intra-frame matching, inter-frame matching, and temporally consistent label assignment. These steps are explained next and visualized in Fig. 3.

Intra-frame matching: At each training iteration, we first generate a large set of candidate masks that can be transferred to neighboring frames. The pseudo-labels from flowIRN might be incomplete, but we expect the predic-

tions from the Mask R-CNN to become more robust as training proceeds. Hence, high-confidence predictions from the model at a given iteration can be used as the candidate set. However, we empirically observed that during early stages of training, masks predicted by Mask R-CNN can be fragmented. We overcome this by also considering the union of all mask predictions which have good overlap with a flowIRN pseudo-label of the same class. The candidate set of predictions after this step includes the original predictions (in practice, we use top 100 predictions) for the frame, as well as the new predictions obtained by combining overlapping predictions, as shown in Fig. 3(a). For a frame at time t , we refer to the original set of predictions from the model as \mathcal{P}^t , and this expanded set as $\mathcal{P}_{\text{exp}}^t$. Each prediction $p_i^t \in \mathcal{P}_{\text{exp}}^t$ corresponds to a triplet: mask, bounding box and the class with highest score for the box, denoted by (m_i^t, b_i^t, c_i^t) respectively.

Inter-frame matching: Next, we wish to transfer some predictions from the current frame t as pseudo-labels to the neighboring frame $t + \delta$ and vice-versa. We only transfer a prediction if it is stably predicted by the current model on both frames. To do this, we first create a bipartite graph between the two frames. The nodes from each frame correspond to the expanded prediction set $\mathcal{P}_{\text{exp}}^t$ and $\mathcal{P}_{\text{exp}}^{t+\delta}$ respectively as shown in Fig. 3(b). The edge weight $e_{ij}^{t,t+\delta}$ between prediction p_i^t and $p_j^{t+\delta}$ is defined as:

$$e_{ij}^{t,t+\delta} = \mathbb{I}(c_i^t = c_j^{t+\delta}) \cdot \text{IoU}(W_{t \rightarrow t+\delta}(p_i^t), p_j^{t+\delta}),$$

where $W_{t \rightarrow t+\delta}$ is a bi-linear warping function that warps the prediction from one frame to another based on the optical flow between them (explained in the appendix). The edge weight is non-zero only if the two predictions share the same class. The weight is high if the warped mask from frame t has high overlap with the mask in $t + \delta$.

The correspondence between predictions of both frames is then obtained by solving the bipartite graph-matching problem with these edge weights, using the Hungarian algorithm. This results in a one-to-one matching between a subset of predictions from $\mathcal{P}_{\text{exp}}^t$ and $\mathcal{P}_{\text{exp}}^{t+\delta}$. We denote the matching result as $\mathcal{M}^{t,t+\delta} = \{(p_i^t, p_j^{t+\delta})\}$, containing pairs of matched predictions from both frames. This comprises pairs of predictions that are temporally stable.

Temporally consistent labels: We use the predictions from frame t which are matched to some predictions in $t + \delta$ in the previous step to define new pseudo-labels for frame $t + \delta$ as shown in Fig. 3(c). Since there can be a lot of spuriously matched predictions, we only transfer high-quality predictions that have some overlap with the original pseudo-labels in frame t . As the process presented in Alg. 1, let $\mathcal{P}_{\text{IRN}}^t$ be the original set of pseudo-labels obtained from flowIRN for frame t and $\mathcal{M}^{t,t+\delta}$ be the matched prediction pairs between two frames. We transfer only those masks from t which have an overlap greater than 0.5 with any of

the original masks in $\mathcal{P}_{\text{fIRN}}^t$ of the same class. Further, when transferring to $t + \delta$, we (a) warp the mask using optical flow and (b) merge it with the matched prediction in $t + \delta$ as shown in Fig. 3(c) to ensure that the mask is not partially transferred. This new set of labels transferred from t to $t + \delta$ are denoted by $\mathcal{P}_{\text{maskCon}}^{t+\delta}$. The steps are explained below. Here, $Merge(\cdot)$ simply takes the union of masks from two predictions to form a new prediction.

Algorithm 1: Temporally consistent assignment

```

Input:  $\mathcal{M}^{t,t+\delta}$ ,  $\mathcal{P}_{\text{fIRN}}^t$ 
Output:  $\mathcal{P}_{\text{maskCon}}^{t+\delta}$ 
1  $\mathcal{P}_{\text{maskCon}}^{t+\delta} \leftarrow \{\}$ 
2 for  $(p_i^t, p_j^{t+\delta}) \in \mathcal{M}^{t,t+\delta}$  do
3   for  $p_{\text{fIRN}}^t \in \mathcal{P}_{\text{fIRN}}^t$  do
4     if  $\text{IoU}(b_i^t, b_{\text{fIRN}}^t) > 0.5$ ,  $c_i^t = c_{\text{fIRN}}^t$  then
5        $p_m^{t+\delta} \leftarrow \text{Merge}(\text{W}_{t \rightarrow t+\delta}(p_i^t), p_j^{t+\delta})$ 
6        $\mathcal{P}_{\text{maskCon}}^{t+\delta} \leftarrow \mathcal{P}_{\text{maskCon}}^{t+\delta} \cup \{p_m^{t+\delta}\}$ 
7       break
8   end
9 end
10 return  $\mathcal{P}_{\text{maskCon}}^{t+\delta}$ 

```

Simultaneously, new pseudo-labels $\mathcal{P}_{\text{maskCon}}^t$ are obtained for t by transferring predictions from $t + \delta$ in a similar fashion. We combine them with the original pseudo-labels from flowIRN to obtain the final set of pseudo-labels $\mathcal{P}_{\text{maskCon}} \cup \mathcal{P}_{\text{fIRN}}$. We also note that while combining these two sets of labels, it is important to suppress smaller masks that are contained within the others. Concretely, we apply non-maximal suppression (NMS) based on an Intersection over Minimum (IoM) threshold. IoM is calculated between two masks as the intersection area over the area of the smaller mask. This avoids label redundancy and helps improve performance as we demonstrate later in ablation experiments. The merged pseudo-labels are used as supervision to train the Mask R-CNN model as shown in Fig. 2, without altering the Mask R-CNN in any way.

Our overall MaskConsist approach does not require extra forward or backward pass, and only adds a small overhead to the original Mask R-CNN during training. During inference, the matching is unnecessary and MaskConsist works similar to Mask R-CNN.

4. Experiments

Unless otherwise specified, models in this section are trained only with frame-level class labels and do not use bounding-box or segmentation labels. We evaluate our model on two tasks: frame-level instance segmentation and video-level instance segmentation. We report performances on two popular video datasets.

4.1. Datasets

Youtube-VIS (YTVIS) [58] is a recently proposed benchmark for the task of video instance segmentation. It

contains 2,238 training, 302 validation, and 343 test videos collected from YouTube, containing 40 categories. Every 5th frame in the training split is annotated with instance segmentation mask.

As the annotation of validation and test splits are not released and only video-level instance segmentation performance is available on the evaluation server, we hold out a subset of videos from the original training split by randomly selecting 10 videos from each category. This results in a train_val split of 390 videos (there are videos belonging to multiple object categories) to conduct frame-level and video-level instance segmentation evaluations. The remaining 1,848 videos are used as the train_train split.

Cityscapes [12] contains high-quality pixel-level annotations for 5,000 frames collected in street scenes from 50 different cities. 19 object categories are annotated with semantic segmentation masks and 8 of them are annotated with instance segmentation masks. The standard 2,975 training frames and their neighboring $t - 3$ and $t + 3$ frames are used for training, and the 500 frames in validation split are used for evaluation.

4.2. Implementation details

Optical flow network: We use the self-supervised DDFlow [33] for optical flow extraction. The model is pre-trained on “Flying Chairs” dataset [14] and then fine-tuned on YTVIS or Cityscapes training videos in an unsupervised way. The total training time is 120 hours on four P100 GPUs and the average inference time per frame is 280ms.

flowIRN: To get flow-amplified CAMs, we set the amplification co-efficient $A = 2$ and threshold $T = \text{Percentile}_{0.8}(\|\mathcal{F}(x)\|_2)$ for YTVIS, and $A = 5$ and $T = \text{Percentile}_{0.5}(\|\mathcal{F}(x)\|_2)$ for Cityscapes. The optical flow is extracted between two consecutive frames (frame t and $t + 1$). The regularization weight λ is set to 2. We train the network for 6 epochs. Other training and inference hyper-parameters are set the same as in [6]. Empirically, we observe that IRN (and flowIRN) is limited by lack of good CAMs when trained only on Cityscapes data. Hence, for experiments on Cityscapes, we train the first-stage of all weakly supervised models (before the Mask R-CNN/MaskConsist stage) first on PASCAL VOC 2012 [15] training-split and then fine-tune on Cityscapes.

MaskConsist: We use ResNet-50 as the backbone, initialized with ImageNet pre-trained weights. For both datasets, the bounding-box IoU threshold is set at 0.5 for intra-frame matching, and IoM-NMS threshold at 0.5 for label combining. The model is trained for 90K iterations for YTVIS, and 75K iterations for Cityscapes, with base learning rate $lr = 0.002$. SGD optimizer is used with step schedule $\gamma = 0.1$, decay at 75% and 88% of total steps. The temporal consistency is calculated between frame t and $t + 5$ ($\delta = 5$) for YTVIS, frame t and $t + 3$ ($\delta = 3$) for Cityscapes. Inference on one frame (short side 480px) takes

Methods	Video Info	Supervision	AP_{50}
Mask R-CNN [17]	✗	Mask	78.24
WSIS-BBTP [20]	✗	Bbox	46.80
WISE [27]	✗	Class	24.54
F2F [29]+MCG [41]	✓	Class	26.31
IRN [6]	✗	Class	29.64
IRN [6]+F2F[29]	✓	Class	30.27
Ours	✓	Class	34.66
Ours (self-training)	✓	Class	36.00

Table 1. Frame-level instance segmentation performance (AP_{50}) on YTVIS train_val split.

210ms. Nvidia Tesla P100 GPU is used in training and test. All hyper-parameters for flowIRN and MaskConsist are selected based on the performance on a small held-out validation split of the corresponding training set.

Experiment setup: On YTVIS, all methods are trained using the training frames (every 5th frame) in train_train split. On Cityscapes, all methods are trained with training frames (frame t) and their two neighboring frames ($t-3$ and $t+3$). Unless otherwise specified, our model is trained in two-steps: first train flowIRN on training frames, then use the pseudo-labels generated by the flowIRN on the training frames to train MaskConsist. For fair comparison, all baseline methods are also trained in two steps: first train the weakly supervised model (e.g., IRN) with frame-level class labels, then use pseudo-labels obtained to train a Mask R-CNN model. This is common practice in weakly supervised segmentation works [6, 27], and improves AP_{50} of all models by at least 2% in our experiments. The same hyper-parameters reported in the original work or published code are retained for all baselines.

We also observe that a three-step training process, where the masks generated by our MaskConsist model are used to train another MaskConsist model, further improves performance. We refer to this as *ours self-training*. Note that unlike other baselines, this involves an additional round of training. On other baseline methods, we also attempted self-training: another round of training using pseudo-labels from the trained Mask R-CNN. However, this either degraded or did not improve performance on the validation set.

During frame-level inference, the trained MaskConsist or Mask R-CNN (for other baselines) is applied on each frame with score threshold of 0.05 and NMS threshold of 0.5 to obtain prediction masks. For video-level evaluation, we apply an unsupervised tracking method [37] on per-frame instance mask predictions to obtain instance mask tracks, with the same hyper-parameters as the original work. We will release our code after paper acceptance.

4.3. Frame instance segmentation

First, we compare frame-level performance with existing instance segmentation models on YTVIS and Cityscapes.

Evaluation metrics: On both YTVIS and Cityscapes, the average precision with mask intersection over union (IoU) threshold at 0.5 (AP_{50}) is used as the metric for in-

Methods	Supervision	Instance seg	Semantic seg
Mask R-CNN [17]	Mask	38.73	79.23
WISE [27]	Class	10.51	35.82
F2F [29]+MCG [41]	Class	10.73	33.26
IRN [6]	Class	12.33	33.48
IRN [6]+F2F[29]	Class	12.53	34.17
Ours	Class	16.05	39.88
Ours (self-training)	Class	16.82	41.31

Table 2. Frame-level instance segmentation (AP_{50}) and semantic segmentation (IoU) on Cityscapes validation split.

stance segmentation. Cityscapes is a popular benchmark for semantic segmentation and we also report the semantic segmentation performance using standard IoU metric.

Baselines: To the best of our knowledge, there is no existing weakly supervised instance segmentation model designed for videos. Existing works are designed for still images and report results on standard image benchmarks like [15]. To compare with these models on video data, we train them (where code is available) with independent video frames of YTVIS or Cityscapes. We also extend existing weakly supervised “video” semantic segmentation models to perform instance segmentation. For upper-bound comparisons, we report results from Mask R-CNN [17] trained with ground truth masks, and WSIS-BBTP [20] trained with bounding box annotations. We list the baselines below and more details can be found in the appendix:

- *WISE [27]*: train on independent frame with class label.
- *IRN [6]*: train on independent frame with class label.
- *F2F [29] + MCG [41]*: use videos with class labels to train F2F to obtain semantic segmentation and combine MCG proposals to obtain instance-level masks as in [61].
- *F2F [29] + IRN [6]*: use optical flow to aggregate CAMs as in F2F to train IRN.

Results: Results on YTVIS are shown in Tab. 1. All methods use two-step training as stated in the experiment setup. WISE and F2F+MCG both use processed CAMs as weak labels and combine results with object proposals (MCG) to distinguish instances. Comparing WISE and F2F+MCG, F2F uses video information that boosts its performance by around 1.8%. IRN+F2F is the closest comparison to our approach, since it is also built on top of IRN and uses video information. Our model outperforms IRN+F2F by more than 4%, and can also benefit from an additional round of self-training (Ours self-training). However, we do not observe any gains when training the Mask R-CNN for another round for other methods.

In Tab. 2, we report frame-level instance segmentation and semantic segmentation results on Cityscapes. For instance segmentation, our method outperforms WISE and IRN by more than 3.7% under AP_{50} . We convert the instance segmentation results to semantic segmentation by merging instance masks of the same class and assigning labels based on scores. On semantic segmentation, our method still outperforms IRN by a large margin.

Methods		Train_Val Split					Validation Split				
		mAP	AP_{50}	AP_{75}	AR_1	AR_{10}	mAP	AP_{50}	AP_{75}	AR_1	AR_{10}
Fully supervised learning methods	IoUTracker+ [58]	-	-	-	-	-	23.6	39.2	25.5	26.2	30.9
	DeepSORT [57]	-	-	-	-	-	26.1	42.9	26.1	27.8	31.3
	MaskTrack [58]	-	-	-	-	-	30.3	51.1	32.6	31.0	35.5
Weakly supervised learning methods	WISE [27]	8.7	22.1	5.5	9.8	10.7	6.3	17.5	3.5	7.1	7.8
	IRN [6]	10.8	26.4	7.7	12.6	14.4	7.3	18.0	3.0	9.0	10.7
	Ours	14.1	34.4	9.4	16.0	17.9	10.5	27.2	6.2	12.3	13.6

Table 3. Video instance segmentation results on Youtube-VIS dataset.

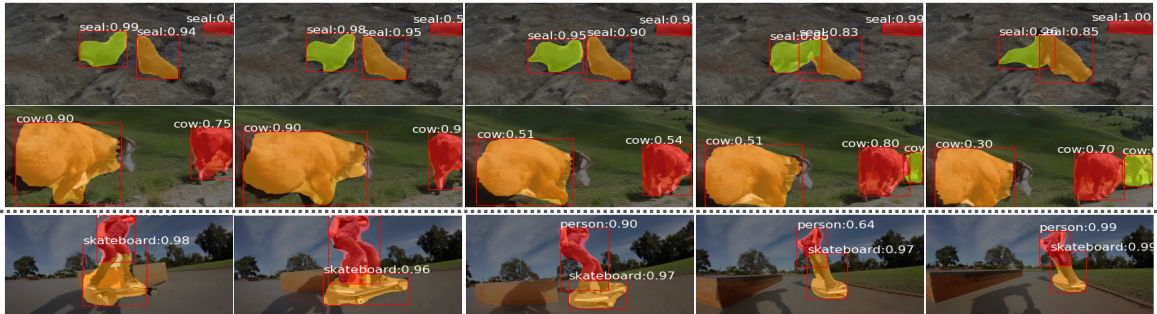


Figure 4. Example Video instance segmentation results from our method on Youtube-VIS dataset.

4.4. Video instance segmentation

Given per-frame instance segmentation predictions, we apply the Forest Path Cutting algorithm [37] to obtain a mask-track for each instance and report VIS results.

Evaluation metric: We use the same metrics as [58]: mean average precision for IoU between [0.5, 0.9] (mAP), average precision with IoU threshold at 0.5 / 0.75 (AP_{50}/AP_{75}), and average recall for top 1 / 10 (AR_1 / AR_{10}). As each instance in a video contains a sequence of masks, the computation of IoU uses the sum of intersections over the sum of unions across all frames in a video. The evaluation is carried out on YTVIS train_val split using YTVIS code (<https://github.com/youtubevos>), and also on YTVIS validation split using the official YTVIS server.

Baselines: Since there is no existing work on weakly supervised video instance segmentation, we construct our own baselines by combining the tracking algorithm in [37] with two weakly supervised instance segmentation baselines: WISE [27] and IRN [6]. We also present published results from fully supervised methods [58, 57] for reference.

As presented in Tab. 3, our model outperforms IRN and WISE by a large margin. On the AP_{50} metric, there is a boost of more than 8% on both train_val and validation splits. We also observe that the performance gap between WISE and IRN decreases compared with frame-level results in Tab. 1, implying temporal consistency is important to realize gains in video instance segmentation. Note that the fully supervised methods are first trained on MSCOCO [32] and then fine-tuned on YTVIS training split, while ours is only trained on YTVIS data. Qualitative VIS results from our method are shown in Fig. 4. Our method generates temporally stable instance predictions and is able to capture different overlapping instances. One failure case

	YTVIS	Cityscapes
IRN [6]	25.42	8.46
IRN+f-Bound	26.60	9.51
IRN+f-CAMs	27.47	10.55
flowIRN	28.45	10.75

Table 4. Ablation study of flowIRN components. Results are reported on training data to evaluate pseudo-label quality. No second-step Mask R-CNN or MaskConsist training is applied here. is shown in the bottom row. As *skateboard* and *person* always appear and move together in YTVIS, our assumption on different instances having different motion is not valid. Thus, these two instances are not well distinguished.

4.5. Effect of modeling temporal information

Our framework explicitly models temporal information in both flowIRN and MaskConsist modules. We explore the effectiveness of each module in this section.

Ablation study of flowIRN: In Tab. 4, we present the instance segmentation results (AP_{50}) of different flowIRN variants. All models are directly tested on the training data to evaluate pseudo-label quality and no second-step training is used in this experiment. Compared to original IRN[6], both flow-amplified CAMs (f-CAMs) and flow-boundary loss (f-Bound) incorporate optical flow information and improve IRN performance. Combining the two leads to our design of flowIRN, which improves by 3.03% on YTVIS and 2.29% on Cityscapes compared to IRN.

In Fig. 5, we show two qualitative examples of incorporating f-CAMs and f-Bound. In the left example, the car (in the circle) moves fast and is partially missed by IRN. After applying f-CAMs, the whole object is well captured in the segmentation mask. In the second example (right column), IRN fails to separate two overlapping persons while the boundary is recognizable in optical flow. After applying f-Bound loss, two instances are correctly predicted.

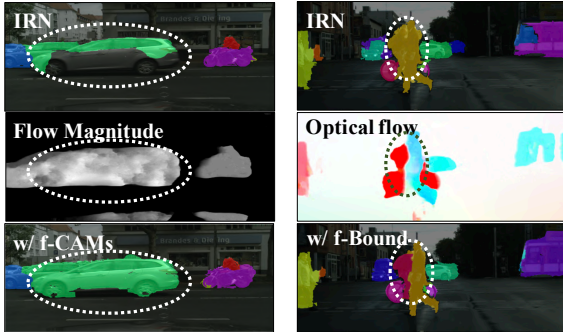


Figure 5. Improvement introduced by f-CAMs and f-Bound. Top: output of IRN. Middle: optical flow extracted for the input frame. Bottom: output after incorporating f-CAMs or f-Bound.

MaskConsist Components			AP_{50}	
Intra-F	Inter-F	IoM-NMS	YTVIS	Cityscapes
✗	✗	✗	31.43	14.66
✗	✓	✓	33.75	14.92
✓	✗	✓	31.08	14.43
✓	✓	✗	33.65	15.27
✓	✓	✓	34.66	16.05

Table 5. Ablation study of MaskConsist components. The numbers in this table are generated by models with two-step training.

Ablation study of MaskConsist: In Tab. 5, we explore the contribution of different components of MaskConsist by disabling one of the three components each time. We observe that inter-frame matching plays the most important role in MaskConsist. It enables the model to incorporate temporal consistency during training and achieves the largest performance boost. IoM-NMS helps avoid false positives corresponding to partial masks from inter-frame matching and improves the performance on top of intra-frame and inter-frame matching. Our best results on both datasets are achieved by combining all three components.

In Tab. 6, we further explore the effectiveness of MaskConsist module by combining it with other weakly supervised instance segmentation methods: WISE [27] and IRN [6]. Cross in the “w/ MC” column denotes the use of Mask R-CNN instead of MaskConsist. The results show that, by incorporating mask matching and consistency in the second stage of training, MaskConsist module consistently improves original weakly supervised methods by about 2%. Combining flowIRN module with MaskConsist achieves the best performance on both YTVIS and Cityscapes.

We also quantitatively evaluate how consistent the predictions of MaskConsist are on consecutive frames. As presented in the fifth column of Tab. 6, we report the temporal consistency (TC) metric similar to [35]. This metric measures the AP_{50} between mask predictions and flow warped masks on consecutive frames in YTVIS. We observe consistent improvement in TC by adding MaskConsist to training.

In Fig. 6, we present two examples of Mask R-CNN and MaskConsist predictions on YTVIS clips. Both models are trained with flowIRN pseudo-labels. Mask R-CNN predic-

Methods	w/ MC	AP_{50}		TC
		YTVIS	Cityscapes	
WISE [27]	✗	24.54	10.51	72.08
	✓	27.03	12.26	76.27
IRN [6]	✗	29.64	12.33	80.98
	✓	31.51	14.72	82.04
Ours	✗	31.43	14.66	80.43
	✓	34.66	16.05	84.36

Table 6. MaskConsist works on top of different weakly supervised instance segmentation methods and improves both AP_{50} and TC .

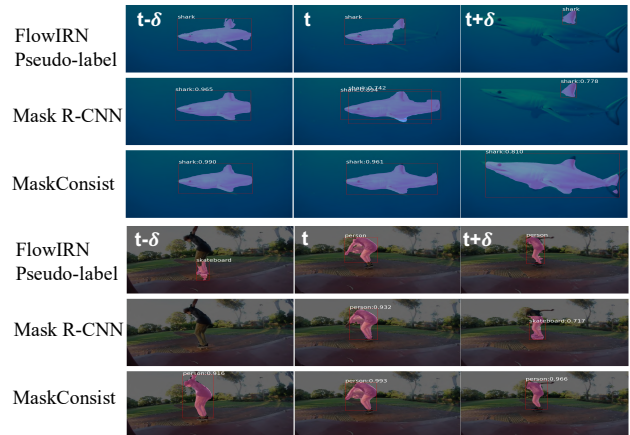


Figure 6. Comparison of Mask R-CNN and MaskConsist on YTVIS. Both models are trained from flowIRN pseudo-label.

tions are more susceptible to noisy pseudo-labels and less consistent across frames, while MaskConsist achieves more stable segmentation results.

Further discussion: Regarding the two types of errors presented in Fig. 1, we observe that our model has larger relative improvement over IRN on more strict metric: 27.3% (12.06% vs. 9.47%) on AP_{75} , compared with 16.9% (34.66% vs. 29.64%) on AP_{50} , indicating our model generates more accurate mask for high IoU metric. While our method outperforms IRN on AP_{50} , our method also predicts more instances per frame (avg. 1.81 instances for ours vs. avg. 1.50 instances for IRN), indicating our method is able to predict more instances with higher accuracy. These demonstrate that the two problems of partial segmentation and missing instance are both alleviated in our model.

5. Conclusion

We observed that image-centric weakly supervised instance segmentation methods often segment an object instance partially or miss an instance completely. We proposed the use of temporal consistency between frames in a video to address these issues when training models from video frames. Our model (a) leveraged the constraint that pixels from the same instance move together, and (b) transferred temporally stable predictions to neighboring frames as pseudo-labels. We proved the efficacy of these two approaches through comprehensive experiments on two video datasets of different scenes. Our model outperformed the state-of-the-art approaches on both datasets.

Appendix A. Supplementary Material

A.1. Explanation of using optical flow gradient

Here we explain that spatial gradient of optical flow helps identify if two pixels are from the same instance. Such information is encoded in the flow-boundary loss (Eq. (3) in the manuscript).

As shown in Fig. 7, let pixel (x, y) in the image is a projection of point (X, Y, Z) in the physical world with velocity (V_X, V_Y, V_Z) . The 3D motion results in optical flow on the image plane $\mathcal{F}_x, \mathcal{F}_y$. We consider for a short time window, most objects move in parallel to the image plane for Youtube VIS and Cityscapes data thus $V_z \approx 0$. For simple mathematical notation, we consider 3D motion only along X-axis $V_x = 0$. The optical flow along image x axis can be written as:

$$\mathcal{F}_x = \frac{f * V_X}{Z},$$

where f is the camera focal length. We explain why we chose to use difference of optical flow first-order gradient instead of directly using difference of optical flow as followed.

For two neighboring pixels p_i and p_j , if they are from the same rigid object, we have $V_{X_i} = V_{X_j} = V_X$, $d(V_X)/dx = 0$. Then, their optical flow difference is:

$$\mathcal{F}_{x_i} - \mathcal{F}_{x_j} = \frac{f * V_{X_i}}{Z_i} - \frac{f * V_{X_j}}{Z_j} = \left(\frac{1}{Z_i} - \frac{1}{Z_j}\right) f * V_X,$$

which is not equal to zero when there is difference in depth for these two pixels.

We propose to use the difference of spatial gradient of optical flow. The spatial gradient of flow along x axis is defined as:

$$\begin{aligned} \frac{d\mathcal{F}_x}{dx} &= \frac{d(1/Z)}{dx} f * V_X + \frac{d(V_X)}{dx} \frac{f}{Z} \\ &= -\frac{1}{Z^2} \frac{dZ}{dx} f * V_X. \end{aligned}$$

For two neighboring pixels p_i and p_j on the same instance surface, assuming the surface is smooth thus $dZ_i/dx = dZ_j/dx = dZ/dx$, their difference of flow gradient is written as:

$$\frac{d\mathcal{F}_{x_i}}{dx} - \frac{d\mathcal{F}_{x_j}}{dx} = -\left(\frac{1}{Z_i^2} - \frac{1}{Z_j^2}\right) \frac{dZ}{dx} f * V_X$$

In practice, the two pixels are in local neighbor and depth values Z_i, Z_j are often large, thus $1/Z_i^2 - 1/Z_j^2 \approx 0$ is better approximation than $1/Z_i - 1/Z_j \approx 0$. This indicates using difference of optical flow gradient is a better signal for pixel affinity calculation.

A similar inference can also be applied to y axis. In practice, we calculate the first-order gradient difference on both directions and encourage the norm to be zero.

A.2. Prediction warping with optical flow

In order to warp Mask R-CNN predictions with optical flow, we first warp the predicted masks using bi-linear interpolation as used in Spatial Transformer Network [23]. The warped mask is then converted to binary mask at threshold of 0.5. Then the bounding box of warped prediction is generated from the warped mask and class label is directly copied. In practice, we implement the mask warping function using `torch.nn.functional.grid_sample` in PyTorch [40] framework.

A.3. Training details for baseline methods

WISE [27]: We train WISE using code published in [3] and [5]. For YTVIS, the model is trained for 20 epochs with learning rate starting at 0.01. For Cityscapes, the model is pre-trained on PASCAL VOC 2012 and then fine-tuned for 40 epochs with learning rate starting at 0.001. The MCG proposals are generated using code in [2]. We generate pseudo-labels on the training split and train Mask R-CNN model as stated in Sec. 4.2 Implementation Details in the manuscript.

F2F [29]+MCG [41]: We first use F2F to generate semantic segmentation masks and then combine with MCG to generate instance masks as pseudo-labels. To train F2F on YTVIS and Cityscapes, we follow the F2F paper to generate the aggregated CAMs: warp the CAMs from 5 consecutive frames to the key frame. The weakly supervised network backbone used in F2F

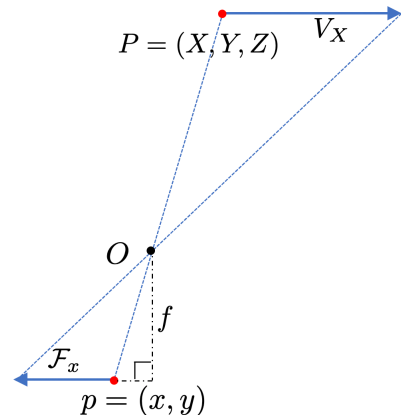


Figure 7. 3D motion results in 2D optical flow following imaging principle.

is not available and we use SEC [26] (code in [4]) which is advised by the F2F authors. For YTVIS, the model is trained for 10,000 iterations with learning rate starting at 0.001. For Cityscapes, the model is pre-trained on PASCAL VOC 2012 and then fine-tuned for 5,000 epochs with learning rate starting at 0.001. To combine with MCG proposals, we adopt a similar approaches in WISE[27]. The resulting instance masks are used as pseudo-labels to train Mask R-CNN.

IRN [6]: We train IRN using code published in [1]. For YTVIS, the model is trained for 6 epochs with learning rate starting at 0.1. For Cityscapes, the model is pre-trained on PASCAL VOC 2012 and then fine-tuned for 6 epochs with learning rate starting at 0.01. The other training and inference parameters are set as default. The resulting instance masks on the training data are then used as pseudo-labels to train Mask R-CNN.

IRN [6]+F2F [29]: We use the flow-warped CAMs as in F2F to train IRN: warp CAMs from 5 neighboring frames to key frame to generate aggregated CAMs. The aggregated CAMs are then used to replace the CAMS used in original IRN. The training parameters are set the same as training original IRN model.

A.4. More visualization of video intance segmentation results

More video instance segmentation results are presented in Fig. 8. Every 5th frame in a video clip from YTVIS train_val split are presented in each row. In the top nine examples, our method generates consistent instance masks with good object coverage and accurate instance boundary. We also include two failure cases in the last two rows, where the object is heavily occluded or object classes have strong co-occurrence pattern.



Figure 8. More video instance segmentation results from our proposed method.

References

- [1] Implementation of IRN. <https://github.com/jiwoon-ahn/irn>. 10
- [2] Implementation of MCG. <https://github.com/jponttuset/mcg>. 9
- [3] Implementation of PRM. <https://github.com/chuchienshu/ultra-thin-PRM>. 9
- [4] Implementation of SEC. https://github.com/halbielee/SEC_pytorch. 10
- [5] Implementation of WISE. https://github.com/ElementAI/wise_ils. 9
- [6] Jiwoon Ahn, Sunghyun Cho, and Suha Kwak. Weakly supervised learning of instance segmentation with inter-pixel relations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2209–2218, 2019. 1, 2, 3, 5, 6, 7, 8, 10
- [7] Aditya Arun, CV Jawahar, and M Pawan Kumar. Weakly supervised instance segmentation by learning annotation consistent instances. *arXiv preprint arXiv:2007.09397*, 2020. 1, 2
- [8] Ali Athar, Sabarinath Mahadevan, Aljoša Ošep, Laura Leal-Taixé, and Bastian Leibe. Stem-seg: Spatio-temporal embeddings for instance segmentation in videos. *arXiv preprint arXiv:2003.08429*, 2020. 2
- [9] Gedas Bertasius and Lorenzo Torresani. Classifying, segmenting, and tracking object instances in video with mask propagation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2
- [10] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiao-xiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4974–4983, 2019. 1
- [11] Hisham Cholakkal, Guolei Sun, Fahad Shahbaz Khan, and Ling Shao. Object counting and instance segmentation with image-level supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12397–12405, 2019. 1, 2
- [12] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2, 5
- [13] Jifeng Dai, Kaiming He, and Jian Sun. Boxesup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1635–1643, 2015. 2
- [14] A. Dosovitskiy, P. Fischer, E. Ilg, P. Häusser, C. Hazırbaş, V. Golkov, P. v.d. Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2015. 5
- [15] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 5, 6
- [16] Weifeng Ge, Sheng Guo, Weilin Huang, and Matthew R Scott. Label-penet: Sequential label propagation and enhancement networks for weakly supervised instance segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3345–3354, 2019. 1, 2
- [17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2961–2969, 2017. 1, 2, 6
- [18] Seunghoon Hong, Donghun Yeo, Suha Kwak, Honglak Lee, and Bohyung Han. Weakly supervised semantic segmentation using web-crawled videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7322–7330, 2017. 2
- [19] Qibin Hou, PengTao Jiang, Yunchao Wei, and Ming-Ming Cheng. Self-erasing network for integral object attention. In *Advances in Neural Information Processing Systems*, pages 549–559, 2018. 2
- [20] Cheng-Chun Hsu, Kuang-Jui Hsu, Chung-Chi Tsai, Yen-Yu Lin, and Yung-Yu Chuang. Weakly supervised instance segmentation using the bounding box tightness prior. In *Neural Information Processing Systems*, 2019. 2, 6
- [21] Anthony Hu, Alex Kendall, and Roberto Cipolla. Learning a spatio-temporal embedding for video instance segmentation. *arXiv preprint arXiv:1912.08969*, 2019. 2
- [22] Zilong Huang, Xinggang Wang, Jiasi Wang, Wenyu Liu, and Jingdong Wang. Weakly-supervised semantic segmentation network with deep seeded region growing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7014–7023, 2018. 2
- [23] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, pages 2017–2025, 2015. 9
- [24] Bin Jin, Maria V Ortiz Segovia, and Sabine Susstrunk. Webly supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3626–3635, 2017. 2
- [25] Anna Khoreva, Rodrigo Benenson, Jan Hosang, Matthias Hein, and Bernt Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 876–885, 2017. 2
- [26] Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *European Conference on Computer Vision*, pages 695–711, 2016. 1, 2, 3, 10
- [27] Issam H Laradji, David Vazquez, and Mark Schmidt. Where are the masks: Instance segmentation with image-level supervision. *The British Machine Vision Conference*, 2019. 1, 2, 6, 7, 8, 9, 10
- [28] Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In

- Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 5267–5276, 2019. 2
- [29] Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Frame-to-frame aggregation of active regions in web videos for weakly supervised semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6808–6818, 2019. 2, 6, 9, 10
- [30] Justin Liang, Namdar Homayounfar, Wei-Chiu Ma, Yuwen Xiong, Rui Hu, and Raquel Urtasun. Polytransform: Deep polygon transformer for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9131–9140, 2020. 2
- [31] Di Lin, Jifeng Dai, Jiaya Jia, Kaiming He, and Jian Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3159–3167, 2016. 2
- [32] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 7
- [33] Pengpeng Liu, Irwin King, Michael R Lyu, and Jia Xu. Ddflow: Learning optical flow with unlabeled data distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8770–8777, 2019. 5
- [34] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8759–8768, 2018. 1
- [35] Yifan Liu, Chunhua Shen, Changqian Yu, and Jingdong Wang. Efficient semantic video segmentation with per-frame inference. *ECCV*, 2020. 8
- [36] Xiankai Lu, Wenguan Wang, Jianbing Shen, Yu-Wing Tai, David J Crandall, and Steven CH Hoi. Learning video object segmentation from unlabeled videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8960–8970, 2020. 2
- [37] Jonathon Luiten, Idil Esen Zulfikar, and Bastian Leibe. Unovost: Unsupervised offline video object segmentation and tracking. In *Proceedings of the IEEE Winter Conference on Applications in Computer Vision*, 2020. 2, 6, 7
- [38] Eslam Mohamed, Mahmoud Ewaisha, Mennatullah Siam, Hazem Rashed, Senthil Yogamani, and Ahmad El-Sallab. Instancemotseg: Real-time instance motion segmentation for autonomous driving. *arXiv preprint arXiv:2008.07008*, 2020. 2
- [39] George Papandreou, Liang-Chieh Chen, Kevin P Murphy, and Alan L Yuille. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1742–1750, 2015. 2
- [40] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 9
- [41] Jordi Pont-Tuset, Pablo Arbelaez, Jonathan T Barron, Ferran Marques, and Jitendra Malik. Multiscale combinatorial grouping for image segmentation and object proposal generation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(1):128–140, 2016. 2, 6, 9
- [42] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. ” grabcut” interactive foreground extraction using iterated graph cuts. *ACM transactions on graphics (TOG)*, 23(3):309–314, 2004. 2
- [43] Fatemeh Sadat Saleh, Mohammad Sadegh Aliakbarian, Mathieu Salzmann, Lars Petersson, and Jose M Alvarez. Bringing background into the foreground: Making all classes equal in weakly-supervised video semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2125–2135, 2017. 2
- [44] Yunhang Shen, Rongrong Ji, Yan Wang, Yongjian Wu, and Liujuan Cao. Cyclic guidance for weakly supervised joint detection and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 697–707, 2019. 1, 2
- [45] Wataru Shimoda and Keiji Yanai. Self-supervised difference detection for weakly-supervised semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5208–5217, 2019. 2
- [46] Chunfeng Song, Yan Huang, Wanli Ouyang, and Liang Wang. Box-driven class-wise region masking and filling rate guided loss for weakly supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3136–3145, 2019. 2
- [47] Guolei Sun, Wenguan Wang, Jifeng Dai, and Luc Van Gool. Mining cross-image semantics for weakly supervised semantic segmentation. In *European Conference on Computer Vision*, 2020. 2
- [48] Meng Tang, Abdelaziz Djelouah, Federico Perazzi, Yuri Boykov, and Christopher Schroers. Normalized cut loss for weakly-supervised cnn segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1818–1827, 2018. 2
- [49] Pavel Tokmakov, Karteek Alahari, and Cordelia Schmid. Learning semantic segmentation with weakly-annotated videos. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [50] Yi-Hsuan Tsai, Guangyu Zhong, and Ming-Hsuan Yang. Semantic co-segmentation in videos. In *European Conference on Computer Vision*, pages 760–775, 2016. 2
- [51] Paul Vernaza and Manmohan Chandraker. Learning random-walk label propagation for weakly-supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7158–7166, 2017. 2
- [52] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking emerges by col-

- orizing videos. In *European Conference on Computer Vision*, pages 391–408, 2018. 2
- [53] Xiaolong Wang, Allan Jabri, and Alexei A Efros. Learning correspondence from the cycle-consistency of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2566–2576, 2019. 2
- [54] Xiang Wang, Huimin Ma, and Shaodi You. Deep clustering for weakly-supervised semantic segmentation in autonomous driving scenes. *Neurocomputing*, 381:20–28, 2020. 2
- [55] Yunchao Wei, Jiashi Feng, Xiaodan Liang, Ming-Ming Cheng, Yao Zhao, and Shuicheng Yan. Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1568–1576, 2017. 1, 3
- [56] Yunchao Wei, Huaxin Xiao, Honghui Shi, Zequn Jie, Jiashi Feng, and Thomas S Huang. Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7268–7277, 2018. 2
- [57] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *Proceedings of IEEE International Conference on Image Processing*, pages 3645–3649, 2017. 7
- [58] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5188–5197, 2019. 2, 5, 7
- [59] Xiaolin Zhang, Yunchao Wei, Jiashi Feng, Yi Yang, and Thomas S Huang. Adversarial complementary learning for weakly supervised object localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1325–1334, 2018. 1, 3
- [60] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. 1, 2
- [61] Yanzhao Zhou, Yi Zhu, Qixiang Ye, Qiang Qiu, and Jianbin Jiao. Weakly supervised instance segmentation using class peak response. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3791–3800, 2018. 1, 2, 6
- [62] Yi Zhu, Yanzhao Zhou, Huijuan Xu, Qixiang Ye, David Doremann, and Jianbin Jiao. Learning instance activation maps for weakly supervised instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3116–3125, 2019. 1, 2