# Self-supervised Training of Proposal-based Segmentation via Background Prediction

**Isinsu Katircioglu**[1], **Helge Rhodin**[1], **Victor Constantin**[1], **Jörg Spörri**[2],
**Mathieu Salzmann**[1], **and Pascal Fua**[1]
`isinsu.katircioglu@epfl.ch`

[1]Computer Vision Laboratory
EPFL
Lausanne, Switzerland

[2]Sports Medical Research Group
Balgrist University Hospital, University of Zurich
Zurich, Switzerland

## Abstract

While supervised object detection methods achieve impressive accuracy, they generalize poorly to images whose appearance significantly differs from the data they have been trained on. To address this in scenarios where annotating data is prohibitively expensive, we introduce a self-supervised approach to object detection and segmentation, able to work with monocular images captured with a moving camera. At the heart of our approach lies the observation that segmentation and background reconstruction are linked tasks, and the idea that, because we observe a structured scene, background regions can be re-synthesized from their surroundings, whereas regions depicting the object cannot. We therefore encode this intuition as a self-supervised loss function that we exploit to train a proposal-based segmentation network. To account for the discrete nature of object proposals, we develop a Monte Carlo-based training strategy that allows us to explore the large space of object proposals. Our experiments demonstrate that our approach yields accurate detections and segmentations in images that visually depart from those of standard benchmarks, outperforming existing self-supervised methods and approaching weakly supervised ones that exploit large annotated datasets.

## 1 Introduction

Recent object detection and segmentation methods have reached impressive precision and recall rates when trained and tested on large annotated datasets [17]. However, large and varied datasets do not warrant the best possible performance in a particular application domain, as each comes with its own challenges and opportunities. While domain-specific models are thus needed, it is impractical to annotate separate and sufficiently large datasets for deep learning.

Therefore, weakly- and self-supervised detection and segmentation of salient foreground objects in complex scenes has recently gained attention [6, 8, 5, 21]. These methods promise effortless processing of community videos with little human intervention. However, a closer look at existing techniques reveals that they make strong assumptions ranging from target objects being on top of static background or relying on pre-trained object localization, object-boundary detection, and optical flow networks. This severely limits their applicability in practice.

To develop a more generic technique, we start from the observation that in most images the background forms a consistent, natural scene. Therefore the appearance of any background patch can be predicted from its surroundings. By contrast, a salient object's appearance is unpredictable from the neighboring scene content and can be expected to be very different from what an inpainting algorithm would

Preprint. Under review.

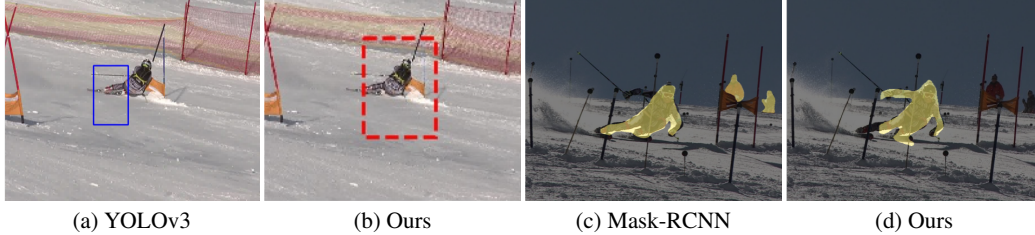|           |          |             |          |
|-----------|----------|-------------|----------|
| (a) YOLOv3 | (b) Ours | (c) Mask-RCNN | (d) Ours |

Figure 1: **Domain specific detection examples.** Our self-supervised method detects the skier well, while YOLO trained on a general dataset does not generalize to this challenging domain. We also compare to MaskRCNN, which succeeds on the skier but detects false positives.

produce. We incorporate this insight into a proposal-generating deep network whose architecture is inspired by those of YOLO [20] and MaskRCNN [10] but does not require explicit supervision.

For each proposal, we synthesize a background image by masking out the corresponding region and inpainting it from the remaining image. The loss function we minimize favors the largest possible distance between this reconstructed background and the input image. This encourages the network to select regions that cannot be explained from their surrounding and are therefore salient. To handle the discrete nature of the proposals, we introduce a Monte Carlo-based strategy to train our network. It operates on a discrete distribution, is unbiased, exhibits low variance, and is end-to-end trainable.

We demonstrate the effectiveness of our unsupervised method on several datasets captured with increasingly mobile cameras, ranging from static to pan-tilt-zoom and hand-held. We will show that our approach applies to images acquired in conditions significantly more general than those of standard benchmarks, without requiring *any* manual annotation. Thus, as shown in Fig. 1, it approaches the quality and sometimes outperforms state-of-the-art detectors that have been trained on large annotated datasets. Retraining or fine-tuning these methods on this data could be done but would require supervision that is hard to obtain, which makes a self-supervised approach attractive. We will make our code and **Handheld190k** dataset publicly available upon acceptance of the paper.

## 2 Related Work

Most salient object detection and segmentation algorithms are fully-supervised [20, 10, 24, 3] and require large annotated datasets with paired images and labels. Our goal is a purely self-supervised method that succeed without segmentation and object bounding box annotations. This is different from those methods requiring domain-specific annotation at training but not at test time, which are often also referred to as *unsupervised object detection* methods [11]. We focus our discussion on neural methods and refer to [15] for the discussion of methods using hand-crafted optimization.

**Weakly-supervised methods.** A classical weakly-supervised example is the Hough Matching algorithm [4]. It uses an object classification dataset and identifies foreground as the image regions that have re-occurring Hough features within images of the same class. Similar principles have been followed using deep networks trained for object detection [28, 14], optical flow [27], and object saliency [16]. These methods make the implicit assumption that the background varies across examples and can therefore be excluded as noise. This assumption is violated in the targeted case of training on domain-specific images, where foreground and background are similar across examples.

**Motion-based methods** Given video sequences, the temporal information can be exploited by assuming that the background changes slowly [2] and linearly [25]. However, even a static scene induces non-homogeneous deformations under camera translation, and it can be difficult to handle all types of camera motion (pan, tilt, zoom) at different pastes, cuts within a single video, and distinguish articulated human motion from background motion [23]. By iteratively refining the crude background subtraction results from [25] with an assemble of student and teacher networks [6, 7], some of the resulting errors could be corrected. However, a strong dependency on the teacher used for bootstrapping remains.

Our approach is conceptually closely related to VideoPCA, which models the background as the part of the scene that can be explained by a low-dimensional linear basis [25]. This implicitly assumes that the foreground is harder to model than the background and can therefore be separated as the
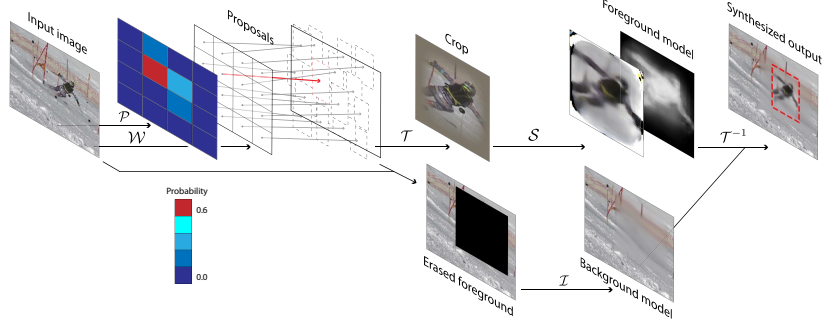
Figure 2: **Method overview.** Encoder-decoder network ($\mathcal{S}$) with an attention mechanism defined by proposal-based detection ($\mathcal{P}, \mathcal{W}$) and spatial transformers ($\mathcal{T}, \mathcal{T}^{-1}$). Carefully designed objective functions make it possible to train this network entirely self-supervised on unknown scenes with a moving background and hand-held camera via an inpainting network ($\mathcal{I}$).

non-linear residual. Instead of motion cues, we introduce a new assumption on the predictability of image patches from the spatial neighborhood.

**Self-supervised Methods.** Most similar to our approach are the self-supervised ones to object-detection [8, 5, 21] that complement auto-encoder networks by an attention mechanism. These methods first detect one or several bounding boxes whose content is extracted using a spatial transformer [13]. This content is then passed through an auto-encoder and re-composited with a background. In [21] the background is assumed to be static and in [8, 5] even single colored, a severe restriction in practice. Crawford et al. [5] use a proposal-based network similar to ours, but resort to approximating the proposal distribution with a continuous one to be differentiable. We demonstrate that much simpler importance sampling is sufficient. By contrast to all of these methods, our approach works with images acquired using a moving camera and given an arbitrarily colored background.

In addition to object detection, the algorithm of [21] also returns instance segmentation masks by reasoning about the extent and depth ordering of multiple people in a multi-camera scene. However, this requires multiple static cameras and a static background at training time, as does the approach of [1] that performs instance segmentation in crowded scenes.

## 3 Method

Our goal is to learn a salient object detector and segmentor from unlabeled videos acquired in as generic a setup as possible, including using hand-held cameras. At inference time, our algorithm takes a single image $\mathbf{I} \in \mathbb{R}^{W \times H}$ as input and outputs a bounding box $\mathbf{b}_m \in \mathbb{R}^{2 \times 2}$, in terms of center location, width, and height, and segmentation mask $\mathbf{S} \in \mathbb{R}^{128 \times 128}$ within that window. Internally, it is a multi-stage process, as visualized in Fig. 2. A first network predicts a set of $C$ candidate object locations $(\mathbf{b}_c)_{c=1}^{C}$ and corresponding probabilities $(p_c)_{c=1}^{C}$. We use a fully-convolutional architecture similar to YOLO [20], as used by the supervised state-of-the-art methods. Second, the $\mathbf{b}_m$ with highest probability $p_m \geq p_c$, $\forall c$ is chosen and its content is decoded into foreground $\hat{\mathbf{I}} \in \mathbb{R}^{128 \times 128}$, segmentation mask $\mathbf{S}$, and background $\mathbf{B} \in \mathbb{R}^{W \times H}$ with separate auto-encoder branches.

### 3.1 Self-Supervised Training

Given a set of unlabeled images $(\mathbf{I}_i)_{i=1}^{N}$, our goal is to train two neural networks $\mathcal{W}(\mathbf{I})$ and $\mathcal{P}(\mathbf{I})$ to propose suitable bounding box candidates and, respectively, the probabilities to select the $\mathbf{b} = \mathcal{W}(\mathbf{I})[c]$ that contains an object. Because object locations are unknown, we do this with an autoencoder objective, $\ell(\mathcal{F}(\mathbf{I}, \mathbf{b}, \mathbf{B}), \mathbf{I})$, that measures how well the autoencoder $\mathcal{F}$ reproduces the input image $\mathbf{I}$ on top of a background $\mathbf{B}$, with the attention on $\mathbf{b}$. As in [5, 21], it is implemented with spatial transformer $\mathcal{T}$ that crops the area of interest, followed by a bottle-neck autoencoder $\mathcal{S}$ that produces foreground and segmentation mask, and second spatial transformer $\mathcal{T}^{-1}$ that undoes the cropping and blends synthesized foreground and background. Formally, we write

$$\mathcal{F}(\mathbf{I}, \mathbf{b}, \mathbf{B}) = \mathcal{T}^{-1}(\hat{\mathbf{I}}, \mathbf{S}, \mathbf{b}, \mathbf{B}), \text{ with } \mathcal{S}(\mathcal{T}(\mathbf{I}, \mathbf{b})) \mapsto (\hat{\mathbf{I}}, \mathbf{S}). \tag{1}$$

3

In the simplest case, $\ell$ is a least-square loss between all pixel values and the background $\mathbf{B}$ is known. Because the attention window $\mathbf{b}$ selects part of the image for decoding, this loss encourages $(\mathcal{W}, \mathcal{P})$ to focus the attention on the object so as to be able to model the foreground on top of $\mathbf{B}$. Moreover, the autoencoder only approximates the actual image, which forces the segmentation mask to just contain those parts not captured by the background.

For now, we assume $\mathbf{B}$ to be known and omit it in our derivations to ease readability, but we will remove this assumption later. In this context, we derive a probabilistic formulation in which not a single, but multiple candidates can give rise to plausible reconstructions. We then reason about the expected loss across all likely candidates,

$$\mathcal{O}(\mathbf{I}) = \mathbf{E}_c \left[ \ell \left( \mathcal{F}(\mathbf{I}, \mathcal{W}(\mathbf{I})[c]), \mathbf{I} \right) \right], \text{ with } c \sim \mathcal{P}(\mathbf{I}), \tag{2}$$

where $\mathbf{E}_c$ denotes the expectation over $c$ drawn from the proposal distribution that is output by the network $\mathcal{P}(\mathbf{I})$. As natural for deep learning, we optimize Eq. 2 with stochastic gradient descent and mini-batches on $(\mathbf{I}_i)_{i=1}^{N}$. Note that minimizing this objective will jointly optimize the detection networks $(\mathcal{W}, \mathcal{P})$ that generate and stochastically select proposals, and the synthesis network $\mathcal{S}$ that models the object appearance and segmentation mask.

Because we have a finite set of candidates, Eq. 2 can be expressed deterministically as a weighted sum over all candidates as

$$\mathcal{O}(\mathbf{I}) = \sum_{c=1}^{C} \mathcal{P}(c|\mathbf{I}) \ell \left( \mathcal{F}(\mathbf{I}, c), \mathbf{I} \right), \tag{3}$$

using the shorthand $\mathcal{F}(\mathbf{I}, c) = \mathcal{F}(\mathbf{I}, \mathcal{W}(\mathbf{I})[c])$. The objective is an explicit function of $\mathcal{P}$ and could be optimized by gradient descent on $\mathcal{O}$. However, this sum is inefficient to evaluate in practice when $C$ is large (e.g., in our experiments $C = 64$, which does not fit in memory). This was also observed by Crawford et al. [5], who resorted to using a continuous approximation of the discrete distribution $\mathcal{P}$ to facilitate end-to-end training. Here, we propose a simpler alternative, exploiting Monte Carlo and importance sampling, which provides an unbiased estimator with low variance.

**Monte Carlo sum.** In principle, the expectation Eq. 2 could be estimated by sampling a small set of $J$ candidate cells from $\mathcal{P}(\mathbf{I})$, for instance, one per mini-batch element, such that

$$\mathcal{O}(\mathbf{I}) \approx \frac{1}{J} \sum_{j=1}^{J} \ell \left( \mathcal{F}(\mathbf{I}, c_j), \mathbf{I} \right), \text{ with } c_j \sim \mathcal{P}(\mathbf{I}). \tag{4}$$

Unfortunately, sampling from such a discrete distribution is not differentiable with respect to its parameters, which precludes end-to-end gradient-based optimization of Eq. 4.

**Importance sampling.** Instead of sampling according to $\mathcal{P}$, we can rewrite Eq. 2 to be over an arbitrary distribution $q$, by reweighting with the quotient of both distributions. That is,

$$\mathcal{O}(\mathbf{I}) = \mathbf{E}_c \left[ \frac{\mathcal{P}(c|I)}{q(c)} \ell \left( \mathcal{F}(\mathbf{I}, c), \mathbf{I} \right) \right], \text{ with } c \sim q. \tag{5}$$

This change of distribution and relative weighting holds for any two probability distributions, as explained in the supplementary material. In practice, we use a mini-batch optimization, with a single sample drawn from $q$ per image, i.e., $J = 1$, for estimating

$$\mathcal{O}(\mathbf{I}) \approx \frac{1}{J} \sum_{j=1}^{J} \left( \frac{\mathcal{P}(c_j|I)}{q(c_j)} \ell \left( \mathcal{F}(\mathbf{I}, c_j), \mathbf{I} \right) \right), \text{ with } c_j \sim q. \tag{6}$$

While moving the distribution into the expectation sum provides differentiability, it comes with the drawback of a potentially large variance, i.e., high approximation error for small $J$. For instance, by choosing a uniform sampling distribution $\mathcal{U}$, most of the uniformly drawn samples will have a low probability in $\mathcal{P}$ and therefore negligible influence. To reduce this variance, we leverage importance sampling and set the sampling distribution $q$ to be similar or equivalent to $\mathcal{P}$. Note that the fraction $\frac{\mathcal{P}(c|I)}{q(c)}$ cancels numerically for $q(c) = \mathcal{P}(c|I)$. However, while values cancel, their derivatives do not; the differentiability of $\mathcal{P}$ is maintained. The sampling distribution $q$ must be treated as a constant.

In practice, we select $q$ based on the current estimate of $\mathcal{P}$ to reduce the variance. To prevent division by very small values that could lead to numerical instability, we define the new distribution $q$ as

$$q(c) = \mathcal{P}(c|\mathbf{I})(1 - C\epsilon) + \epsilon \,. \tag{7}$$

As a side effect, $\epsilon$ increases the probability that an unlikely case is chosen, which induces a form of exploration that is helpful in the early training stages of the network. We analyze the attained variance reduction in the supplementary material.

Notably, the gradient of Eq. 6 equals that of the likelihood ratio method [9] used in the REINFORCE algorithm [29]. In reinforcement learning terms [26], setting $q = \mathcal{P}$ would correspond to a single step of on-policy learning. We refrained from motivating our derivation in this manner because the simple importance sampling rule is sufficient to explain our approach.

## 3.2 Training with Dynamic Cameras

Having derived an efficient training scheme for proposal-based segmentation when $\mathbf{B}$ is given, we would like to reduce the more difficult moving camera scenario to the former. This requires predicting the background image, which, in the absence of prior shape and appearance information, requires to identify and inpaint the pixels that are different from the background.

While this task entails object detection, our primary goal, the related, yet simpler inpainting task can easily be trained by removing a region and predicting it from its immediate surrounding [18, 30]. A network, $\mathcal{I}$, trained on this self-supervised inpainting task, would nonetheless not reconstruct foreground objects if fully removed in the input because the surrounding background gives no cues of their presence. Detecting foreground objects can therefore be cast as finding the area $\mathbf{b}$ that, when inpainted, yields the largest image reconstruction error.

To accomplish this search efficiently, we define the background objective probabilistically, similar to the foreground we write,

$$\mathcal{G}(\mathbf{I}) = -\mathbf{E}_c \left[ \ell(\mathcal{I}(\mathbf{I}, \mathbf{b}), \mathbf{I}) \right], \text{ with } c \sim \mathcal{P}(\mathbf{I}), \tag{8}$$

where $\mathcal{I}$ takes the image $\mathbf{I}$ and the region $\mathbf{b}$ to inpaint as input and $\ell, \mathcal{P}$, and $\mathbf{b} = \mathcal{W}[c]$ are as before. Note that opposed to the foreground objective, we use the negative expectation. This negative reconstruction error encourages the selection of those regions were the true image is dissimilar to the reconstructed background when minimizing Eq. 8.

However, as illustrated in Figure 4(d), straight regression with a neural network or brute force search would yield trivial solutions. For instance, erasing extensively large regions, containing an object or not, will lead to higher inpainting errors, just because of the increased number of reconstructed pixels. To prevent this, we reformulate Eq 8 so as to compute differences only within the cropped region and normalize the result by the crop area. However, this, in turn, favors locations with high error density, irrespectively of their size, as shown in Figure 4(b).

To overcome degenerate cases, we combine the new background objective $\mathcal{G}$ with the foreground objective $\mathcal{O}$ of Eq. 2, substituting the known $\mathbf{B}$ with learned inpainting $\mathcal{I}(\mathbf{I}, \mathbf{b})$. The reason behind this is that these two terms are complementary: While $\mathcal{G}$ prefers locations that cover the object neither precisely nor entirely, $\mathcal{O}$ favors a tight fit over partial coverage but has a trivial solution when $\mathbf{b}$ is on a background region, i.e., not covering the object and having nothing to encode.

We exploit the complementary behavior of these two terms by partitioning their influence on the individual network components. For $\mathcal{G}$, we limit the gradient flow to only update $\mathcal{P}$, freezing the remaining network modules $\mathcal{W}$ and $\mathcal{S}$. By contrast, we use $\mathcal{O}$ to update only $\mathcal{W}$ and $\mathcal{S}$. Thereby, Eq. 8 controls the coarse localization through detection, while Eq. 2 provides fine-grained regression of $\mathbf{b}$. Note that this separation into coarse and fine localization is only possible with the chosen proposal-based detection framework; direct regression to bounding boxes, as in [21], would preclude this important separation.

**Inpainting network.** In principle, any off-the-shelf inpainting network trained on large and generic background datasets could be used. For instance, [18, 30] can produce very plausible results. However, they hallucinate objects and are therefore ill-suited to our goal. Instead, we train $\mathcal{W}$ from scratch, by reconstructing rectangular, randomly removed image regions, as shown in Figure 2. This network will

|        |        |        |        |
|:------:|:------:|:------:|:------:|
|  (a)   |  (b)   |  (c)   |  (d)   |

Figure 3: **Off-the-shelf inpainting results,** on skiing. (a) Input image with the hidden middle part, followed by inpainting with (b) [18], (c) [30] trained on ImageNet. (d) and [30] trained on Places2.

attempt to memorize all the images in the training set. Nevertheless, as for generic inpainting, moving objects that are independent of the surroundings's cannot be reconstructed. Note that overfitting of this network to the training set is acceptable, if not intended, as it is not needed at inference time.

**Implementation details.** Since no labels for intermediate supervision are available, we found naive end-to-end training to be unreliable. To counteract this, we use ImageNet-trained weights for initialization; rely on a perceptual loss on top of the per pixel $\ell_2$ loss for $\ell$; exploit the Focal Spatial Transformer (FST) of [21] to speed up convergence; and scale the erased region in $\mathcal{I}$ to be 1.1 of that predicted by $\mathcal{W}$ to increase the chances of covering the object. Moreover, we limit the location offsets to 1.5 the cell width and discard those outside the image. In addition, we rely on $\ell_2$ priors on the output of $\mathcal{P}$ and $\mathcal{W}$, and an $\ell_1$ prior on $\mathbf{S}$. The pixel reconstruction and perceptual losses are weighted 1:2, and the priors have a weight of 0.1, 1, and 0.1, respectively, to compensate for their different magnitudes. Additional details are given in the supplemental document.

## 4   Experiments

In this section, we evaluate our approach to self-supervised salient object detection and segmentation. Note that our algorithm works on single images at inference time and only requires the background inpainting model at training time. Even though our approach is not people specific, we focus here on people-detection because this is the only domain for which there are benchmark datasets that contain relatively long sequences of the same scene, which is what our method requires for training purposes. First, we use a controlled environment with arbitrary but static background to compare our method to a state-of-the-art self-supervised one and to perform an ablation study. Then, we use skiing footage acquired using PTZ-cameras along with footage of people performing 14 everyday activities recorded using hand-held cameras to demonstrate that existing supervised methods that do well in the controlled environment struggle to adapt to such challenging conditions, whereas our approach delivers promising results. We provide additional results in the supplementary material.

### 4.1   People in a Controlled Environment

We compare our method against state-of-the-art ones on the **Human3.6m** dataset [12] that comprises 3.6 million frames and a set of 15 motion classes. It features nine subjects, five for training and two for validation, seen from different viewpoints against a static background and with good illumination.

**Comparative Results.** On the left side of Table 1, we compare our detection accuracy to that of a very recent self-supervised deep learning method [21], using the mean detection precision (mAP), the mean precision of having an intersection-over-union (IoU) of more than 50%. Our slightly lower accuracy stems from not explicitly assuming a static background, which [21] does. It is valid in a lab but results in total failure in outdoor scenes with moving backgrounds, such as those discussed below.

**Ablation Study.** Here we show that our model choices for training and probabilistic inference are important. Using uniform sampling instead of importance sampling, as described in Eq. 6, does not converge, as shown in Fig. 4(a). Fig. 4(b) shows that joint instead of our separated training of $\mathcal{P}$ and $\mathcal{W}$ with $\mathcal{O}$ and $\mathcal{G}$ produces bounding boxes that are too large. Fig. 4(c) shows that using only the background objective leads to small detections that miss the subject and (d) that direct regression without multiple candidates diverges. These failure cases are representative for the whole dataset.

**Multiple people.** Although our focus is on handling single objects or persons, our probabilistic framework can handle several at test time by sampling more than once. Fig. 5 shows the predicted cell probability as blue blobs whose size is proportional to it. The fully-convolutional architecture
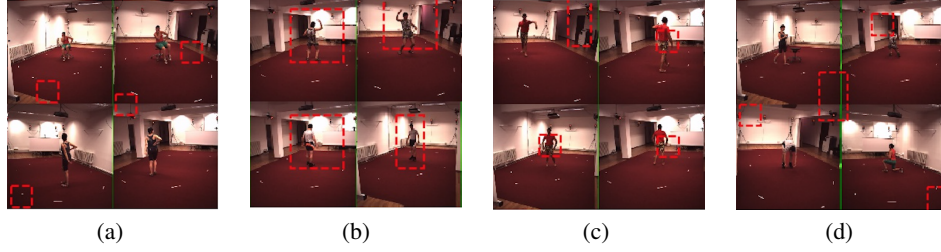
Figure 4: **Ablation study on H36M.** (a) Uniform sampling does not converge. (b) Joint training of $\mathcal{O}$ and $\mathcal{G}$ (c) only $\mathcal{G}$ (d) direct regression of a single bounding box using $\mathcal{O}$ and $\mathcal{G}$.



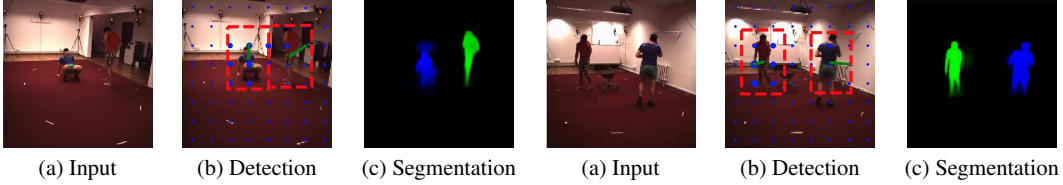| (a) Input | (b) Detection | (c) Segmentation | (a) Input | (b) Detection | (c) Segmentation |

Figure 5: **Multi-person detection and segmentation results**, generated by sampling our model multiple times. As the model is trained on single persons this only works for non-intersecting cases.

operates locally and thereby predicts a high person probability next to both subjects. As a result, both the detection and segmentation results remain accurate so long persons are sufficient separated.

## 4.2 Skiers Filmed Using a PTZ-Cameras

We now turn to the out-of-the-ordinary motions of six skiers on a slalom course featured in the **Ski-PTZ-Dataset** dataset [22]. The six skiers are split four/one/one to form training, validation, and test sets; totaling to, respectively, 7800, 1818 and 1908 frames. The intrinsic and extrinsic parameters of the pan-tilt-zoom cameras are constantly adjusted to follow the skier. As a result, nothing is static in the images, the background changes quickly, and there are additional people standing in the back.

We use the full image as input, evaluate detection accuracy in relation to the available 2D pose annotation, and segmentation accuracy by manually segmenting 36 frames (one from each of the six cameras and two test sequences). As shown on the right side of Table 1, our method delivers a mAP$_{0.5}$ score that is significantly better than that of the general YOLO [20] detector trained on MSCOCO. For an analysis of the segmentation quality, we ran several related works on this dataset



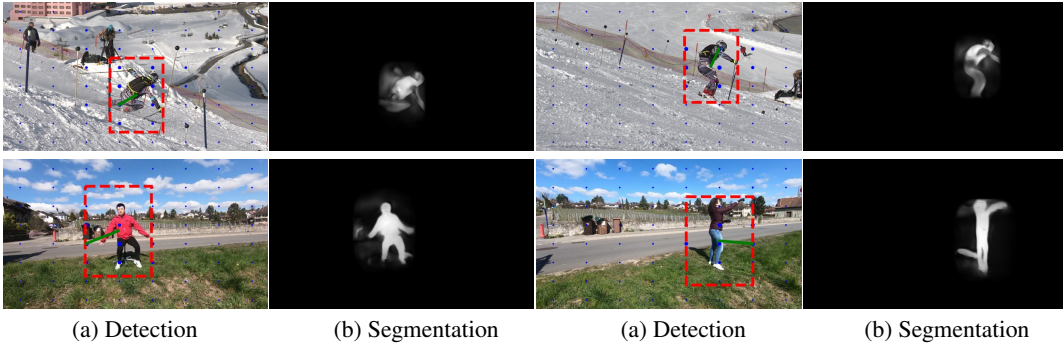| (a) Detection | (b) Segmentation | (a) Detection | (b) Segmentation |

Figure 6: **Qualitative results on Ski-PTZ-camera and Handheld190k.** Example results on the test images. (a) The detection results show the predicted bounding box with red dashed lines, the relative confidence of the grid cells with blue dots and the bounding box center offset with green lines. (b) Soft segmentation mask predictions. Note that in the second row, the moving clouds are not segmented but the shadow of the person can be included.

| H36M dataset | | Ski-PTZ-Dataset | |
| --- | --- | --- | --- |
| Method | $mAP_{0.5}$ | Method | $mAP_{0.5}$ |
| NSD [21] | **0.710** | YOLOv3 [20] | 0.155 |
| Ours | 0.580 | Ours | **0.278** |

Table 1: **Detection** results on the **H36M** and **Ski-PTZ-Dataset** datasets. They are expressed in terms of $mAP_{0.5}$, the mean probability of having an intersection-over-union (IOU) of more than 50%.

| | Ski-PTZ-Dataset | | | | Handheld190k | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Method | Precision | Recall | F Measure | J Measure | Precision | Recall | F measure | J measure |
| MaskRCNN [10] [d] | 0.75 | 0.65 | 0.68 | 0.65 | **0.91** | **0.82** | **0.86** | **0.77** |
| ARP [15][w] | **0.94** | **0.76** | **0.83** | **0.73** | 0.87 | 0.64 | 0.73 | 0.63 |
| VideoPCA [25][s] | 0.49 | **0.84** | 0.61 | 0.56 | 0.33 | **0.91** | 0.48 | 0.49 |
| Unsup-DilateU-Net [7][s] | 0.74 | 0.76 | **0.74** | **0.65** | **0.83** | 0.75 | **0.79** | **0.70** |
| Ours[s] | **0.75** | 0.56 | 0.63 | 0.56 | 0.75 | 0.74 | 0.74 | 0.66 |

Table 2: **Segmentation** results on the **Ski-PTZ-Dataset** and **Handheld190k** dataset. Ours exceeds or is on par with the self-supervised methods (marked with [s]), and approaches the accuracy of weakly-supervised (marked with [w]) and fully-supervised methods (marked with [d]).

and list them in Table 2, in terms of precision, recall, F-, and J-measure as defined in [19]. To be fair, we compensate for different segmentation masks quantification levels by a grid search (at 0.05 intervals) to select the best threshold in terms of J-measure for each method. This measure is defined as the intersection-over-union between the ground-truth segmentation mask and the prediction. The F-Score is the harmonic average between the precision and the recall on the mask boundaries.

Interestingly, MaskRCNN trained on a large generic dataset is outperformed by ARP on this dataset. without using any object localization data, our method is on par with MaskRCNN and close to weakly supervised methods that train on large datasets with motion boundary and segmentation mask annotation, while exceeding the self-supervised one of [25]. Our results exceed the only existing self-supervised object segmentation method using deep learning [7] in precision but are slightly behind in recall, F- and J-measures. Part of this difference can be attributed to [7] using a segmentation mask discriminator that is trained on the combination of the ImageNet VID and YouTube Objects datasets. Albeit also trained in a self-supervised fashion, it thereby leverages additional information and results are not one-to-one comparable.

Further qualitative results are shown in Figure 6. The probability distribution, visualized as blue dots that increase in magnitude with the predicted likelihood, show clear peaks on the persons. Limitations are the slightly blurred and bleeding masks and occasional false positives, reducing precision.

### 4.3 Daily Activities Captured Using Hand-Held Cameras

We introduce a new **Handheld190k** dataset that features three training and two validation sequences that comprise $120\,000$ and $69\,000$ images, respectively, with a single actor performing actions mimicking those introduced in **H36M**. The camera operators moved laterally, to test robustness to camera translation and hand-held rotation. We provide examples of our detection and segmentation results in Fig. 6, more are given in the supplemental document. Our method is robust to the undirected camera motion and to dynamic background motion, such as branches swinging in the wind and clouds moving at this windy day, and to salient textures in the background, such as that of the house facade.

To perform a quantitative comparison, we manually segmented 36 validation images taken from six different motion classes—-pose, phone, shopping, directions, petting a dog, and greeting—-with the subject in many different poses. We then ran several existing methods on this dataset and evaluated the same quantities as for skiing. In this scenario, MaskRCNN yields the highest scores, which is not surprising as the tested sequences are similar to its training set. It is closely followed by [7], which however uses a discriminator that is trained unsupervised on another, larger dataset. Compared to skiing, we exceed ARP [15] in F- and J-measure and have a larger margin on VideoPCA [25], both often fail in separating the non-homogeneously moving background due to the hand-held camera motion.

# 5 Conclusion

We have proposed a self-supervised method for object detection and segmentation that lends itself for application in domains where general purpose detectors fail. Our core contributions are the Monte Carlo-based optimization of proposal-based detection, new foreground and background objectives, and their joint training on unlabeled videos captured by static, rotating and handheld cameras. Our latest experiments demonstrate that, even if trained only on single persons, our approach generalizes to multi-person detection, as long as the persons are sufficiently separated. In contrast to many existing solutions [2, 23, 6], our approach does not exploit temporal cues. In the future, we will integrate temporal dependencies explicitly, which will facilitate addressing the scenario where multiple people interact closely, by incorporating physics-inspired constraints enforcing plausible motion.

# References

[1] P. Baqué, F. Fleuret, and P. Fua. Deep Occlusion Reasoning for Multi-Camera Multi-Target Detection. In *International Conference on Computer Vision*, 2017.

[2] O. Barnich and M. Van Droogenbroeck. Vibe: A universal background subtraction algorithm for video sequences. *IEEE Transactions on Image processing*, 20(6):1709–1724, 2011.

[3] J. Cheng, Y.-H. Tsai, S. Wang, and M.-H. Yang. Segflow: Joint learning for video object segmentation and optical flow. In *Proceedings of the IEEE international conference on computer vision*, pages 686–695, 2017.

[4] M. Cho, S. Kwak, C. Schmid, and J. Ponce. Unsupervised object discovery and localization in the wild: Part-based matching with bottom-up region proposals. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1201–1210, 2015.

[5] E. Crawford and J. Pineau. Spatially invariant unsupervised object detection with convolutional neural networks. In *Conference on Artificial Intelligence*, 2019.

[6] I. Croitoru, S.-V. Bogolin, and M. Leordeanu. Unsupervised learning of foreground object detection. *arXiv preprint arXiv:1808.04593*, 2018.

[7] I. Croitoru, S.-V. Bogolin, and M. Leordeanu. Unsupervised learning of foreground object segmentation. *International Journal of Computer Vision*, pages 1–24, 2019.

[8] S. Eslami, N. Heess, T. Weber, Y. Tassa, D. Szepesvari, K. Kavukcuoglu, and G. Hinton. Attend, infer, repeat: Fast scene understanding with generative models. In *Advances in Neural Information Processing Systems*, pages 3225–3233, 2016.

[9] P. Glynn. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33(10):75–84, 1990.

[10] K. He, G. Gkioxari, P. Dollar, and R. Girshick. Mask R-CNN. In *International Conference on Computer Vision*, 2017.

[11] Y.-T. Hu, J.-B. Huang, and A. G. Schwing. Unsupervised video object segmentation using motion saliency-guided spatio-temporal propagation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 786–802, 2018.

[12] C. Ionescu, I. Papava, V. Olaru, and C. Sminchisescu. Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2014.

[13] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial Transformer Networks. In *Advances in Neural Information Processing Systems*, pages 2017–2025, 2015.

[14] S. D. Jain, B. Xiong, and K. Grauman. Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. In *2017 IEEE conference on computer vision and pattern recognition (CVPR)*, pages 2117–2126. IEEE, 2017.

[15] Y. J. Koh and C.-S. Kim. Primary object segmentation in videos based on region augmentation and reduction. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7417–7425. IEEE, 2017.

[16] S. Li, B. Seybold, A. Vorobyov, A. Fathi, Q. Huang, and C.-C. Jay Kuo. Instance embedding transfer to unsupervised video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6526–6535, 2018.

[17] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Zitnick. Microsoft COCO: Common Objects in Context. In *European Conference on Computer Vision*, pages 740–755, 2014.

[18] D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, and A. Efros. Context Encoders: Feature Learning by Inpainting. *CoRR*, abs/1604.07379, 2016.

[19] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-Hornung, and L. Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017.

[20] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You Only Look Once: Unified, Real-Time Object Detection. In *Conference on Computer Vision and Pattern Recognition*, 2016.

[21] H. Rhodin, V. Constantin, I. Katircioglu, M. Salzmann, and P. Fua. Neural scene decomposition for multi-person motion capture. 2019.

[22] H. Rhodin, J. Spoerri, I. Katircioglu, V. Constantin, F. Meyer, E. Moeller, M. Salzmann, and P. Fua. Learning Monocular 3D Human Pose Estimation from Multi-View Images. In *Conference on Computer Vision and Pattern Recognition*, 2018.

[23] C. Russell, R. Yu, and L. Agapito. Video pop-up: Monocular 3d reconstruction of dynamic scenes. In *European Conference on Computer Vision*, pages 583–598. Springer, 2014.

[24] H. Song, W. Wang, S. Zhao, J. Shen, and K.-M. Lam. Pyramid dilated deeper convlstm for video salient object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 715–731, 2018.

[25] O. Stretcu and M. Leordeanu. Multiple frames matching for object discovery in video. In *BMVC*, volume 1, page 3, 2015.

[26] R. Sutton and A. Barto. *Reinforcement Learning*. MIT Press, 1998.

[27] P. Tokmakov, K. Alahari, and C. Schmid. Learning video object segmentation with visual memory. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4481–4490, 2017.

[28] X.-S. Wei, C.-L. Zhang, J. Wu, C. Shen, and Z.-H. Zhou. Unsupervised object discovery and co-localization by deep descriptor transforming. *arXiv preprint arXiv:1707.06397*, 2017.

[29] R. J. Williams. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. In *Reinforcement Learning*. 1992.

[30] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5505–5514, 2018.