

Fast video object segmentation with Spatio-Temporal GANs

S. Caelles^{1,*} A. Pumarola^{2,*} F. Moreno-Noguer² A. Sanfeliu² L. Van Gool¹
¹Computer Vision Lab, ETH Zürich ²Institut de Robotica i Informàtica Industrial, CSIC-UPC

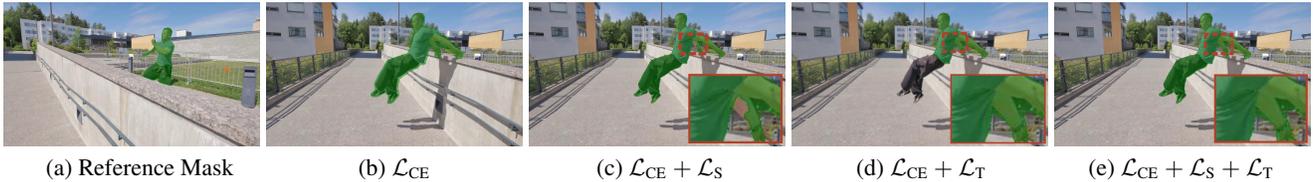


Figure 1: **FaSTGAN Overview.** Our approach learns spatio-temporal object models given a reference mask. When only training with cross-entropy (\mathcal{L}_{CE}) large errors at the object boundary are produced (b). Adding spatial consistency (\mathcal{L}_S) improves the latter but errors still occur on semi-occluded parts (c). Instead, using temporal consistency (\mathcal{L}_T) improves the mask propagation but the model struggles to recover from occluded parts (d). In this work we leverage the benefit of combining spatio-temporal information (e) running at 32 fps.

Abstract

Learning descriptive spatio-temporal object models from data is paramount for the task of semi-supervised video object segmentation. Most existing approaches mainly rely on models that estimate the segmentation mask based on a reference mask at the first frame (aided sometimes by optical flow or the previous mask). These models, however, are prone to fail under rapid appearance changes or occlusions due to their limitations in modelling the temporal component. On the other hand, very recently, other approaches learned long-term features using a convolutional LSTM to leverage the information from all previous video frames. Even though these models achieve better temporal representations, they still have to be fine-tuned for every new video sequence. In this paper, we present an intermediate solution and devise a novel GAN architecture, FaSTGAN, to learn spatio-temporal object models over finite temporal windows. To achieve this, we concentrate all the heavy computational load to the training phase with two critics that enforce spatial and temporal mask consistency over the last K frames. Then at test time, we only use a relatively light regressor, which reduces the inference time considerably. As a result, our approach combines a high resiliency to sudden geometric and photometric object changes with efficiency at test time (no need for fine-tuning nor post-processing). We demonstrate that the accuracy of our method is on par with state-of-the-art techniques on the challenging YouTube-VOS and DAVIS datasets, while running at 32 fps, about $4\times$ faster than the closest competitor.

1. Introduction

The problem of semi-supervised video object segmentation consists in segmenting an object from the background throughout a video sequence given its ground truth mask in the initial frame. Large video datasets like DAVIS [42, 44] and the recently released YouTube-VOS [64] have spurred a number of deep networks methods [4, 41, 56, 38, 9, 28, 22, 61, 50, 65, 64, 60, 7, 11, 24, 2, 5, 32, 35] that improve the performance of approaches from the pre-deep learning era [6, 17, 37, 43, 52] by a large margin. The problem, however, is still far from being solved. Occlusions, rapid object movements, appearance changes and similarity among different instances are still a major obstacle that often require heavy post-processing operations, human intervention and expensive model fine-tuning.

In order to achieve robustness to these challenges, descriptive spatio-temporal object models encoding appearance and geometric changes need to be learned. Most existing state-of-the-art approaches [4, 41, 65] heavily rely on a reference mask to fine-tune the model and in some cases, use the previous mask as a guidance. Formally, if we denote this reference mask by \mathbf{Y}_0 and the RGB frame at time t by \mathbf{X}_t , these approaches model the mask $\hat{\mathbf{Y}}_t$ as $p(\hat{\mathbf{Y}}_t|\mathbf{Y}_0, \mathbf{X}_0, \mathbf{X}_{t-1}, \mathbf{X}_t)$ or $p(\hat{\mathbf{Y}}_t|\mathbf{Y}_0, \mathbf{X}_0, \mathbf{X}_{t-1}, \mathbf{X}_t, \hat{\mathbf{Y}}_{t-1})$. Since no temporal consistency is enforced, these methods tend to be robust to drifting, but they underperform when the object drastically changes its appearance.

This can be remedied by leveraging the temporal consistency of the segmented mask. However, while this was a common practice in the past [17, 37, 43], it is not usual among deep learning methods, in part due to the

*First two authors contributed equally

absence of large scale video object segmentation datasets. Very recently, Xu *et al.* [64] used a convolutional LSTM trained with the YouTube-VOS dataset to learn long-term temporal dependencies from the entire history of the object in the video. That is, the mask \hat{Y}_t is modeled as $p(\hat{Y}_t | Y_0, X_0, X_1, \dots, X_t)$. While this approach demonstrates improved performance compared to previous baselines which did not enforce temporal consistency, it seems to be too generic, as it still needs a computationally demanding fine-tuning step when applied to a sequence with unseen objects.

In this paper, we propose FaSTGAN, an intermediate solution that learns spatio-temporal object appearance models over finite time horizons that does not require fine-tune nor post-processing. Essentially during training, we model the segmentation masks as $p(\hat{Y}_t | Y_0, X_0, X_{t-K}, \dots, X_t)$, where K is the size of the temporal window. In order to implement this model, we design a regressor network architecture inspired by the agile Siamese encoder-decoder structure proposed by Wug *et al.* [60]. In its original form, this regressor is only fed by the reference mask and the masked image at the previous time step. To exploit all information within a temporal window of size K , we could naively make the regressor have access to more information by concatenating features from the K previous frames. This, however, would heavily penalize the efficiency and adaptability of the model. We have therefore devised a novel GAN architecture (Figure 2) in which, during training, this regressor is combined with $K + 1$ discriminators that enforce the temporal and spatial coherence of the generated masks over the temporal window. At test, these discriminators are removed, keeping the original efficiency of the Siamese regressor, while allowing it to model the object across longer time horizons.

As a result, our architecture only uses video data to train and does not require any kind of fine-tuning nor post-processing operations at test time. This makes our approach very efficient, running at 32 fps on 512×512 video frames which is about $4 \times$ faster than [60], which was the fastest video segmentation method so far with 7.7 fps reported in their original work. Furthermore, the accuracy of the segmentation masks we obtain on both DAVIS and YouTube-VOS datasets is on a par with state-of-the-art methods that focus on speed. All code, pre-trained models and pre-computed results used in this paper will be released.

2. Related Work

Video object segmentation: In recent years, video object segmentation has experienced a tremendous increase in popularity due to the publication of large datasets (DAVIS [42, 44] and YouTube-VOS [64]) that have enabled the training of deep learning techniques. The two main set-

tings to tackle this problem are semi-supervised and unsupervised. In the former, the ground truth mask for the object of interest in the first frame of the video sequence is given to the method whereas in the latter no information is given to the algorithm and usually the object with predominant motion is segmented. In this work, we tackle the semi-supervised setting and therefore we focus on previous work in that field.

Traditional approaches used temporal super-pixel [6, 17], optimization in the bilateral space [37], or optimal selection of object proposals [43] to obtain the object segmentation mask for each frame in the video sequence. [4] and [41] were the first two approaches to apply deep learning to the problem. Specifically, [4] fine-tuned the network using the first frame of the video sequence whereas [41] used the mask from the previous frame as an input to the network. [56] extended [4] with an online learning strategy, while [38] also extended [4] by combining its result with Mask-RCNN [20].

Another set of techniques have tried to incorporate the information of the first frame in different ways. [50] formulated the problem as a pattern matching with the initial mask, [65] introduced the initial mask in the network in a batch-norm like layer. [60] used a Siamese network to combine the low level features of the initial frame with the current one together with the back propagation through time training strategy introduced in [23].

Moreover, some methods have tried to leverage metric learning to solve the problem [7, 11, 24], divide the object in multiple parts and track each of them [8], integrate CNN features with traditional energy minimization techniques [2, 5] or design a complex architecture with re-identification and bidirectional propagation modules [32].

Previous approaches mostly rely on single image segmentation, using at most the previous mask as a temporal consistency constraint. There are, however, a few attempts to exploit the temporal dimension better. For instance, some techniques have leveraged optical flow as an additional input to their network. [9] and [28] built a second CNN branch to process optical flow, [22] used it as a prior in the decoder of the network and, [61] used it to align the features from previous frames. While these methods use optical flow priors trained in a separate context, [61] used the large scale YouTube-VOS dataset (released by them) to train a convolutional LSTM in an end-to-end manner.

Our approach lies in between methods that do not use temporal consistency, and [61], which learns long-term dependencies with an RNN.

Generative Adversarial Networks: Since the GAN framework was introduced in [16] to generate synthetic images from noise, its performance has drastically improved in subsequent works [47, 29, 3] achieving almost indistinguishable results from real images. Moreover, it has been suc-

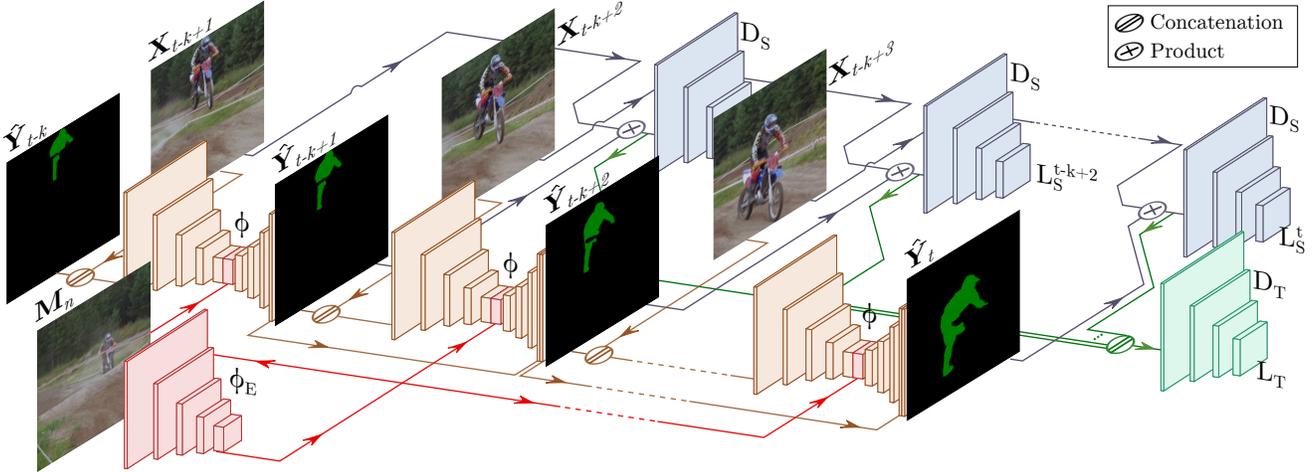


Figure 2: **Overview of our method, FaSTGAN, for video object segmentation.** The diagram displays the model at training time for frames $t - K + 1, t - K + 2, \dots, t$. The proposed architecture consists of three main components: a segmentation regressor ϕ , a spatial critic D_S and a temporal critic D_T . Weights are shared across each network of the same instance. ϕ_E is the encoder part of ϕ and it encodes the image-mask pair $\langle \mathbf{X}_t, \hat{\mathbf{Y}}_{t-1} \rangle$ and \mathbf{M}_n at every time step, see Section 4.1.

successfully applied in a wide range of different tasks such as conditional image synthesis [27, 46], video generation [45, 58], domain adaptation [53], super-resolution [13, 31] or object detection [59].

Recently, several works have tried using GANs for semantic segmentation [34, 51, 36, 25]. [34] trained the discriminator to differentiate between real and predicted probability maps. [36] used two different discriminators to obtain local and global semantic consistency for the human part segmentation problem. [25] tackled the semi-supervised semantic segmentation task using the discriminator to obtain a confidence label map for unlabeled data.

Finally, there have been several applications of GANs in video mainly for video synthesis [57, 54, 48] and conditional video synthesis [58]. In video synthesis, [48] splits the synthesis between two generators, the temporal generator that outputs a latent variable for each frame and the spatial generator that generates the frame from the latent variables. [54] splits the latent space in a motion subspace and a content subspace to generate videos with the same object but performing different motions. Recently, in the conditional setting, [58] synthesizes videos given the dense pose estimation, the semantic maps or the boundaries for each frame in the video sequence.

To the best of our knowledge, we are the first ones to successfully apply GANs to perform *segmentation in videos*. Furthermore, we use a considerably higher image resolution than previous segmentation works (512 vs. 256 in [36] or 321 in [25]) and we leverage recent advances introduced in the images synthesis task, i.e., WGAN with gradient penalty [18], in order to improve stability during training.

3. Problem Formulation

We next formally describe our problem, and generalize the formulation introduced in Section 1 to an arbitrary number of objects, i.e., we aim to design a deep learning model able to segment and track objects along a video sequence given only one single segmentation mask per object.

Let $\mathbf{x} = (\mathbf{X}_1, \dots, \mathbf{X}_T)$ be an input RGB video with T frames, where $\mathbf{X}_t \in \mathbb{R}^{H \times W \times 3}$ denotes the t^{th} frame. Let us also define $\mathbf{m} = (\mathbf{M}_1, \dots, \mathbf{M}_N)$ as a set of reference segmentations of N objects. The reference segmentation $\mathbf{M}_n \in \mathbb{R}^{H \times W \times (3+1)}$ for the n -th object is the concatenation of the first RGB frame in which the object appears with its annotated binary mask. Our goal is to estimate the masks $\hat{\mathbf{y}}$ of all N objects along the entire sequence \mathbf{x} , i.e., we want to learn the mapping $\mathcal{M} : (\mathbf{x}, \mathbf{m}) \rightarrow \hat{\mathbf{y}}$, where $\hat{\mathbf{y}} = (\hat{\mathbf{Y}}_1, \dots, \hat{\mathbf{Y}}_T)$, and $\hat{\mathbf{Y}}_t \in \mathbb{R}^{H \times W \times N}$ contains the N tracked objects masks in the t^{th} video frame. We define the ground truth masks \mathbf{y} for a certain sequence \mathbf{x} as $\mathbf{y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_T)$.

4. Our Approach

Figure 2 shows an overview of FaSTGAN, our proposed approach for video object segmentation. A regressor ϕ is trained on the binary segmentation task of separating the desired object from the background and two WGAN-GP [18] based critics, D_S and D_T , enforce the model to produce semantically and temporally consistent estimates. To simplify the model, we introduce a Markov assumption defining the conditional distribution $p(\hat{\mathbf{y}}|\mathbf{x}, \mathbf{m})$ to be factorized as:

$$p(\hat{\mathbf{y}}|\mathbf{x}, \mathbf{m}) = \prod_{t=1}^T p(\hat{\mathbf{Y}}_t | \mathbf{x}_{t-K}^t, \mathbf{m}, \hat{\mathbf{y}}_{t-K}^{t-1}), \quad (1)$$

meaning that we assume objects in a certain frame to be trackable given the reference segmentations \mathbf{m} , the current and $K-1$ previous frames $\mathbf{x}_{t-K}^t = (\mathbf{X}_{t-K}, \dots, \mathbf{X}_t)$, and the $K-1$ previous estimated segmentation masks $\hat{\mathbf{y}}_{t-K}^{t-1} = (\hat{\mathbf{Y}}_{t-K}, \dots, \hat{\mathbf{Y}}_{t-1})$.

During training, the regressor ϕ is required to learn the mapping \mathcal{M} by modeling the distribution $p(\hat{\mathbf{Y}}_t | \mathbf{x}_{t-K}^t, \mathbf{m}, \hat{\mathbf{y}}_{t-K}^{t-1})$, as $\hat{\mathbf{Y}}_t = \phi(\mathbf{X}_t, \mathbf{m}, \hat{\mathbf{Y}}_{t-1})$. A key property of our design is that the regressor does not directly receive the full information of the K -temporal window. Instead, it is trained to binary segment the K frames, one at a time, given a single foreground/background segmentation \mathbf{M}_n of the desired object $(\mathbf{x}_{t-K}^t, \mathbf{M}_n) \rightarrow \hat{\mathbf{y}}_{t-K}^t$. Each predicted mask is then independently evaluated by a spatial critic $D_S(\mathbf{X}_i, \hat{\mathbf{Y}}_i) \forall i \in [t-K, t]$ that aims to penalize non-consistent semantic masks. The temporal consistency is assessed by a temporal critic $D_T(\mathbf{x}_{t-K}^t, \hat{\mathbf{y}}_{t-K}^t)$ that jointly evaluates the K segmentation masks. Additionally, to feed the regressor with information from previous estimations we introduce a ‘‘temporal skip connection’’ (see details in the following subsection). Note that with this strategy, the regressor is adapted to produce temporally coherent masks within a horizon of size K , without having to simultaneously process the K frames. This will be crucial to deliver a very fast regressor at test time, when the critics will be discarded.

In the following subsections we describe in detail each of these components as well as the proposed training loss.

4.1. Model Architecture

Segmentation Regressor. Given the current frame \mathbf{X}_t , the single-view reference segmentation \mathbf{M}_n of the desired object, and the previous estimate $\hat{\mathbf{Y}}_{t-1}$, the segmentation regressor ϕ aims to separate the desired object from the background producing the current estimate mask $\hat{\mathbf{Y}}_t = \phi(\mathbf{X}_t, \mathbf{M}_n, \hat{\mathbf{Y}}_{t-1})$. We denote the encoder part of ϕ as ϕ_E . Similar to [60], ϕ_E maps the image-mask pairs $\langle \mathbf{X}_t, \hat{\mathbf{Y}}_{t-1} \rangle \in \mathbb{R}^{H \times W \times (3+1)}$ and \mathbf{M}_n to a shared low-dimensional space. Then, feature matching using global convolutions [40] between both features is performed and fed into the decoder part of ϕ to produce the estimated mask $\hat{\mathbf{Y}}_t$. In other words, we train ϕ to refine a rough mask from the previous frame $t-1$ in order to estimate the mask at the current frame t using a reference segmentation of the object \mathbf{M}_n .

In order to enforce temporal consistency along time, we extend the architecture from [60] with a ‘‘temporal skip connection’’. To do so, we concatenate features in the last decoder layer of ϕ with features extracted by the same layer in the previous frame. To reduce the memory complexity involved, we reduce the number of channels in the previous frame feature map by a factor of 1/8 with a 3×3 convolution making the computational cost increase negligible.

Adding this connection not only provides the model with information from previous frames but also acts as a simplified model memory state similar to an RNN. Moreover, when training, the gradients of future frame predictions will directly flow into previous estimates guiding the optimization to take into account that an estimated mask at frame t will have a direct impact into future ones.

Spatial Critic. Partial object segmentation masks and background leaks are two of the most frequent errors in segmentation. To this end, we introduce a spatial critic network, D_S , trained to evaluate the semantic consistency of the pixels in the estimated mask, i.e., we penalize the segmented regions that do not contain one and only one fully covered object ending in its borders without extending into the background. In the experimental section, we prove this supervision to be more informative for the model in comparison to just using the cross entropy loss, as it improves its accuracy. The structure of D_S resembles that of the PatchGan [27] mapping the product¹ $(\mathbf{X}_t \cdot \hat{\mathbf{Y}}_t) \in \mathbb{R}^{H \times W \times 3}$ to an output matrix $\mathbf{S} \in \mathbb{R}^{H/2^6 \times W/2^6}$ where $\mathbf{S}[i, j]$ is used as a partial function to compute the *earth mover’s distance* (EM) between the distributions of the input overlapping patch ij and the real one. This critic helps to improve the difficult task of defining the mask boundaries while enforcing the model not to produce miss-classified small segmentation blobs around the objects of interest.

Temporal Critic. When tracking an object instance across a video sequence we do not only need to have semantically coherent masks, these masks must also be consistent across time. To this end, the *temporal critic* D_T simultaneously evaluates the current estimate w.r.t. to $K-1$ neighbor frames by learning the mapping $(\mathbf{x}_{t-K}^t, \hat{\mathbf{y}}_{t-K}^t) \rightarrow S$, where as above $S \in \mathbb{R}^{H/2^6 \times W/2^6}$ is the overlapping partial scores of a PatchGan based critic. This critic helps to learn relative deformation patterns and plausible absolute motion in the segmentation mask pixel space across time. Also, it enforces the model to generate smooth transitions across mask estimates without large noisy changes.

4.2. Learning the Model

The loss function we define contains three terms, namely a *balanced binary cross entropy* loss to penalize pixel-wise masks errors w.r.t. ground-truth annotations; the *spatial consistency loss* to drive the distribution of the estimates to the distribution of the training masks; and the *temporal consistency loss* that penalizes temporally non-consistent masks.

¹By an abuse of notation we perform the element-wise product on the three RGB channels of \mathbf{X}_t .

Balanced Binary Cross Entropy Loss. We first define the supervised pixel-wise loss for binary classification. To take into account the imbalance between the number of pixels in the object of interest and the background, we apply the balancing strategy proposed in [62] originally used for contour detection. Therefore, the balanced binary cross entropy loss \mathcal{L}_{CE} for K frames is given by:

$$\mathcal{L}_{\text{CE}} = -\frac{1}{K} \sum_{t=K}^t \sum_{j \in \mathbf{Y}_t} \left[\beta \mathbf{Y}_{tj} \log p(\hat{\mathbf{Y}}_{tj} = 1) \right. \\ \left. + (1 - \beta)(1 - \mathbf{Y}_{tj}) \log p(\hat{\mathbf{Y}}_{tj} = 0) \right] \quad (2)$$

where \mathbf{Y}_t is the ground truth binary mask and $\beta = |\mathbf{Y}_t^-|/|\mathbf{Y}_t|$ is the percentage of pixels not belonging to the object.

Spatial Consistency Loss. In order to optimize the *spatial critic* D_S parameters and learn the distribution of the training data, we use the modification of the standard GAN min-max strategy game [16] proposed by WGAN-GP [1]. In our initial experiments, we observed that replacing the Jensen-Shannon (JS) divergence by the continuous Earth Mover Distance resulted in a more stable training. To introduce the required Lipschitz constraint, we apply the gradient penalty proposed by [18] computed as the norm of the gradients with respect to the critic input. Formally, if we denote the data distribution by P_r , the model distribution by P_g , and the random interpolation distribution between masked images by $P_{\tilde{x}}$, the spatial consistency loss \mathcal{L}_S is given by:

$$\mathcal{L}_S = \frac{1}{K} \sum_{t=K}^t \left[\mathbb{E}_{\mathbf{Y}_t \sim P_r} [D_S(\mathbf{X}_t \cdot \mathbf{Y}_t)] - \mathbb{E}_{\hat{\mathbf{Y}}_t \sim P_g} [D_S(\mathbf{X}_t \cdot \hat{\mathbf{Y}}_t)] \right] \\ - \frac{1}{K} \sum_{t=K}^t \lambda_{\text{gp}} \mathbb{E}_{\tilde{x} \sim P_{\tilde{x}}} [(\|\nabla_{\tilde{x}} D_S(\tilde{x})\|_2 - 1)^2], \quad (3)$$

where \tilde{x} is the random interpolation between $\langle \mathbf{X}_t \cdot \mathbf{Y}_t, \mathbf{X}_t \cdot \hat{\mathbf{Y}}_t \rangle$ and λ_{gp} is the penalty coefficient.

Temporal Consistency Loss. With the previously defined losses, the segmentation regressor ϕ is enforced to estimate pixel-wise spatial-consistent masks. However, there is no constraint to guarantee temporal consistency, meaning that the predicted masks should cover similar content across frames. With the *temporal critic* D_T , we push ϕ to maintain temporal consistency by enforcing similarity between joint distributions of K estimated and annotated masks. To estimate the distance between the distributions, we use the approximated Kantorovich-Rubinstein duality [55] of the Earth Mover Distance as proposed in [18]:

$$\mathcal{L}_T = \mathbb{E}_{x \sim P_r} [D_T(x)] - \mathbb{E}_{\hat{x} \sim P_g} [D_T(\hat{x})] \\ - \lambda_{\text{gp}} \mathbb{E}_{\tilde{x} \sim P_{\tilde{x}}} [(\|\nabla_{\tilde{x}} D_T(\tilde{x})\|_2 - 1)^2], \quad (4)$$

where $x = \mathbf{x}_{t-K}^t \cdot \mathbf{y}_{t-K}^t$ and $\hat{x} = \mathbf{x}_{t-K}^t \cdot \hat{\mathbf{y}}_{t-K}^t$ are the real and estimated conditional distributions respectively, and \tilde{x} is the random interpolation between $\langle x, \hat{x} \rangle$. Note that, again, by an abuse of notation, we extended the element-wise product between each \mathbf{x}_t and \mathbf{y}_t (or $\hat{\mathbf{y}}_t$) along the 3 color channels of \mathbf{x}_t .

Overall Loss. To learn to track an object instance across time, we finally define the following minmax problem:

$$\phi^* = \arg \min_{\phi} \max_{D \in \mathcal{D}} (\lambda_{\text{CE}} \mathcal{L}_{\text{CE}} + \lambda_S \mathcal{L}_S + \lambda_T \mathcal{L}_T) \quad (5)$$

where λ_{CE} , λ_S and λ_T are the hyper-parameters that control the relative importance of every loss term and \mathcal{D} the set of 1-Lipschitz functions.

5. Training Details

Our model’s encoder ϕ_E is a ResNet 50 [21] pretrained on ImageNet [12] on the task of image labeling. In order to obtain our final model, we divide the training process in two steps. First, our model is trained only using the supervised loss \mathcal{L}_{CE} to obtain $\phi^* = \min_{\phi} \mathcal{L}_{\text{CE}}$ for 6 epochs on YouTube-VOS [64]. As a result, ϕ^* has an initial understanding of the video object segmentation task and provides a better initialization than ImageNet.

Then, we use ϕ^* as an initialization to train our spatio-temporal model using the loss defined in Eq. 5 in DAVIS17 [44] for 40 epochs. In our experiments, we observe that adding the critics once the model is initialized closer to the final task helps to stabilize training. With the idea to bring the predicted and the ground truth masks distributions in the discriminators closer at each iteration, we overwrite the ground truth pixel values \mathbf{Y}_t with the values of the predicted masks $\hat{\mathbf{Y}}_t$ that are correctly estimated with an uncertainty lower than 0.25. Also, at each iteration, the ground truth masks are augmented by adding Gaussian noise with mean and variance equal to $\hat{\mathbf{Y}}_t$ statistics.

Our model is trained with images of size 512×512 augmented with horizontal flipping, random scaling with factors $[0.75, 1.25]$ and $[-30, 30]$ degrees rotations. Also, at each training iteration, the reference object frame \mathbf{M}_n of a sequence \mathbf{x} is randomly chosen (instead of always being the first frame in which the object appears). We use Adam [30] with a learning rate $1e-5$, β_1 0.5, β_2 0.999, batch size 6 and polynomial decay with power 0.9. During the training of the spatio-temporal model, the learning rate is constant for the first 10 epochs and ϕ is optimized once for every 5 optimization steps of the critic networks. The weight coefficients for the loss terms in Eq. (5) are set to $\lambda_{\text{CE}} = 100$, $\lambda_S = 1$, $\lambda_T = 1$ and $\lambda_{\text{gp}} = 10$.

In order to better approximate the mask error propagation that occurs at test time during training, we set the temporal window size K to the highest value that fits in our

GPU memory, $K = 4$. Note that K is just used during training and only information from the previous frame is used at test time. This parameter is similar to *back propagation through time* introduced in [23] where it was shown that training robustness improves by propagating as many K estimated masks as possible rather than the ground-truth.

We concentrate all computational load to the training stage, which requires 4 NVidia[®] Titan Xp, 3 of them used for training the regressor and 1 for the critics. Our model takes 2 days to finish pretraining on YouTube-VOS and 3 days for the final training on DAVIS17. During test, we only require one single GPU with at least 600Mb of memory. When using an NVidia[®] Titan Xp, we can process videos up to 32 fps.

6. Experimental Evaluation

We thoroughly evaluate our method, FaSTGAN, quantitatively and qualitatively. We compare our approach against current state of the art on semi-supervised video object segmentation: PReMVOS [35], OSVOS^S [38], DyeNet [32], CRN [22], MoNet [61], RGMP [60], MaskRNN [23], VideoMatch [24], YTVOS [63], PML [7], OSMN [65] and BVS [37].

We evaluate our method on the tasks of single object (DAVIS16 [42]) and multiple object video segmentation (DAVIS17 [44], YouTube-VOS [64]). The segmentation accuracy is reported as region similarity (intersection over union \mathcal{J}), contour accuracy (\mathcal{F} measure), and their mean ($\mathcal{J}\&\mathcal{F}$). For the subsets whose annotation is non-public, we compute the results using the submission website provided by the organizers of the challenge.

6.1. Ablation Study

In Table 1, we perform a comprehensive ablation study to analyze the effect of the different loss components that we use in our method during training. In our baseline, we only use the balanced binary cross entropy loss \mathcal{L}_{CE} that penalizes wrong predictions at each pixel and frame independently. Therefore, we do not enforce any spatial or temporal consistency.

First, we introduce the spatial discriminator with its associated loss (\mathcal{L}_S). This improves substantially the accuracy of the method in the contours of the objects boosting the \mathcal{F} measure by more than 1.5 points. This improvement can also be seen qualitatively in Figure 1 comparing (b) to (c).

After that, we replace the previous discriminator by the temporal one with its associated loss (\mathcal{L}_T). As a consequence, performance increases considerably gaining 1.5 points in $\mathcal{J}\&\mathcal{F}$. Now, the model predicts masks with better temporal consistency as can be seen in Figure 1 when comparing the right arm of the man in (c) versus (d). However, the temporal smoothness of the masks enforced by the

		DAVIS16 Val		
\mathcal{L}_S	\mathcal{L}_T	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}
-	-	80.0	79.8	80.2
✓	-	80.5	79.2	81.8
-	✓	81.5	80.7	82.2
✓	✓	81.9	80.2	83.5

Table 1: **Quantitative Ablation Study:** Comparison between the different loss components.

model is sometimes too severe and the regressor has difficulties recovering from large disoccluded parts (Figure 1d).

Finally, we combine the spatial and temporal discriminators and their respective losses. As a result, FaSTGAN is capable of incorporating the accuracy improvements in the contours introduced by the spatial discriminator together with the temporal masks propagation smoothness introduced by the temporal discriminator achieving a final score of 81.9 $\mathcal{J}\&\mathcal{F}$ in DAVIS16. As it can be seen in Figure 1e, the final mask predicted by our spatio-temporal model segments properly the right hand of the man and it can also recover the segmentation of the legs that were occluded in previous frames.

6.2. Evaluation on DAVIS16

Figure 3 shows the accuracy using $\mathcal{J}\&\mathcal{F}$ in DAVIS16 versus the frame rate for various state-of-the-art methods. We can clearly see that our method is in a unique position achieving similar accuracy to previous methods that focused on speed like RGMP, OSMN or PML, while improving on their frame rate by at least a factor of 4. When compared to methods that try to maximize their accuracy (OSVOS^S, MoNet or PReMVOS), we lose around 5 points in $\mathcal{J}\&\mathcal{F}$, but in exchange we increase on their frame rate by a factor of at least 350. As a result, FaSTGAN sets new state of the art in terms of high frame rate while obtaining high accuracy. This brings the field of video object segmentation closer to applications in real time scenarios.

Table 2 reports the results presented in Figure 3 quanti-

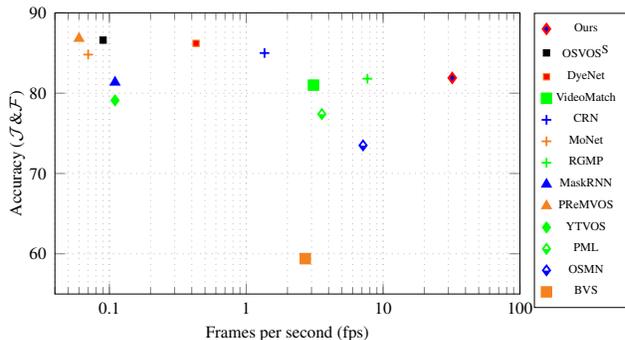


Figure 3: **Accuracy versus speed.** $\mathcal{J}\&\mathcal{F}$ in DAVIS16 with respect to frames per second (fps).

Measure	Ours	RGMP	OSMN	PML	VideoMatch	CRN	DyeNet	YTVOS	MaskRNN	OSVOS ^S	MoNet	PReMVOS	RGMP*
$\mathcal{J}\&\mathcal{F}$ Mean \mathcal{M} \uparrow	81.9	81.8	73.5	77.4	–	85.0	–	–	81.3	86.6	84.8	86.8	79.0
Frames per second \downarrow	32.5[‡] /32.2*/30.3 [†]	7.7*	7.1 [◊]	3.6	3.1 [†]	1.4 [†]	0.43 [‡]	0.11	0.11	0.09 [†]	0.07 [†]	0.06	32.2*
\mathcal{J} Mean \mathcal{M} \uparrow	80.2	81.5	74.0	75.5	81.0	84.4	86.2	79.1	80.4	85.6	84.7	84.9	78.4
Recall \mathcal{O} \uparrow	94.6	91.7	87.6	89.6	–	97.1	–	–	96.0	96.8	96.8	96.1	92.1
Decay \mathcal{D} \downarrow	9.6	10.9	9.0	8.5	–	5.6	–	–	4.4	5.5	6.4	8.8	3.6
\mathcal{F} Mean \mathcal{M} \uparrow	83.5	82.0	72.9	79.3	–	85.7	–	–	82.3	87.5	84.8	88.6	79.7
Recall \mathcal{O} \uparrow	94.3	90.8	84.0	93.4	–	95.2	–	–	93.2	95.9	94.7	94.7	90.8
Decay \mathcal{D} \downarrow	9.1	10.1	10.6	7.8	–	5.2	–	–	8.8	8.2	8.6	9.8	3.6

Table 2: **DAVIS16 Val**: FaSTGAN versus the most recent state of the art, more methods can be found in the DAVIS website². RGMP* is pretrained on YouTube-VOS instead of simulated data. Frames per second reported on a Titan X for [†], Titan Xp for [‡], Quadro M600 for [◊] or 1080Ti for ^{*}, methods without specifier did not report hardware in their publications.

	DAVIS17 Val				YouTube-VOS Val			
	$\mathcal{J}\&\mathcal{F}$	fps	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$	fps	\mathcal{J}	\mathcal{F}
Ours	60.2	16.5[‡] /16.3*/15.4 [†]	57.6	63.4	52.0	17.2[‡] /17.1*/16.1 [†]	50.0	53.9
RGMP	66.7	2.97*	64.8	68.6	52.8	2.66*	51.4	54.2
OSMN	54.8	3.62 [◊]	52.5	57.1	51.8	3.79 [◊]	50.3	52.1
OSVOS ^S	68.0	0.05 [†]	64.7	71.3	–	–	–	–
PReMVOS	77.8	0.03	73.9	81.8	72.2	0.03	69.3	75.2
RGMP*	56.2	16.3*	52.8	59.6	46.3	17.1*	44.3	48.2

Table 3: **DAVIS17 and YouTube-VOS**: FaSTGAN versus the state of the art, more methods can be found in the DAVIS website³. fps stands for frames per second and is computed assuming linear scaling with the number of objects, thus using fps from DAVIS16 and multiplying by the mean number of objects in a certain set. Specifier (*, [†], [‡], [◊], ^{*}) definitions are the same than in Table 2.

tatively and in more depth. In general, methods usually focus either on accuracy or speed at the expense of achieving lower performance in the other. Our method clearly focuses on speed while trying to retain as much accuracy as possible. Compared to the best previous method in the fast speed spectrum, RGMP, we achieve a similar accuracy while improving their frame rate by almost a factor of 4 when tested in the same hardware (32.2 vs 7.7 fps). This improvement is mainly achieved by simplifying their multiscale testing approach by using only a single scale with an image resolution of 512x512. Note that by testing RGMP with our single scale strategy, a similar frame rate is obtained but their accuracy drops by 5 points in $\mathcal{J}\&\mathcal{F}$ which is much lower than our accuracy (76.78 vs 81.9).

In order to show the effect of YouTube-VOS pretraining in previous methods, we train RGMP by substituting their synthetic data generation pretraining with YouTube-VOS video sequences, denoted as RGMP* in Table 2. To do so, we first train their method on YouTube-VOS and then we fine-tuned it on DAVIS17, we stop training in both cases when the loss flattens. In order to provide a fair comparison, our single scale test strategy is used which also improves their fps. Pretraining on YouTube-VOS improves their performance by roughly 2 points in $\mathcal{J}\&\mathcal{F}$ (79.0 vs 76.78) which is still significantly lower than the accuracy of our model.

²https://davischallenge.org/davis2016/soa_compare.html

6.3. Evaluation on DAVIS17 and YouTube-VOS

We also report the performance of FaSTGAN in multi-object video segmentation in Table 3 against state of the art methods in DAVIS17 and YouTube-VOS. Our method is still the fastest at such accuracy outperforming other competitors that are 4 times slower. When compared to RGMP pretrained in YouTube-VOS and tested with our single scale strategy, for the same speed we outperform their model by 4 and 5.7 points in $\mathcal{J}\&\mathcal{F}$ in both DAVIS17 and YouTube-VOS, respectively.

Against the original RGMP model, our method has lower accuracy in DAVIS17, but it has similar results in YouTube-VOS. In DAVIS17, there are several small objects and their multi-scale testing strategy helps to obtain better performance in such a scenario. Moreover, their model in the multi-object scenario not only predicts the masks for objects in the sequence, but also propagates the background as an additional object increasing the average number of objects. As a result of not using multi-scale testing and not propagating an additional object, FaSTGAN runs more than 5 times faster than their method.

Compared to the best performing method, PReMVOS, our inference speed is more than 500 times faster. However, achieving such high frame rate comes at the cost of a drop in accuracy. We believe that optimizing the inference while achieving high accuracy is more challenging in the multi-object scenario compared to the single object one and poses an interesting direction for future works.

6.4. Fairness in method comparison

We would like to briefly discuss the difficulties in providing a fair comparison with other video object segmentation models. First of all, methods in Table 2 and Table 3 have been pretrained in a wide variety of different datasets. For instance, OSVOS^S, DyNet, OSMN, PML and PReMVOS use COCO [33]; MoNet, PReMVOS, DyeNet, VideoMatch, PML, RGMP and CRN use PASCAL VOC [15, 19]; and RGMP uses as well ECSSD [49] and MSRA 10K [10]. Also, MoNet, PReMVOS, DyeNet and CRN use optical

³https://davischallenge.org/davis2017/soa_compare.html



Figure 4: **Qualitative results:** Sample sequences from YouTube-VOS (top 2), DAVIS17 (middle 2) and DAVIS16 (bottom 2). Leftmost image is the initial reference frame, the rest of the images are the predictions for the following frames.

flow produced by Flownet2.0 [26] which is trained using [14, 39]. Before the release of YouTube-VOS, static image datasets were used to train most methods due to the lack of a large scale video object segmentation dataset. We expect future methods in the field to gradually converge to YouTube-VOS pretraining which would make the comparison among different methods easier.

Moreover, previous methods report timings in a wide variety of GPU types. In order to provide a fair comparison, we list the GPU type used in each publication in Table 2 and and Table 3 when available and we test our method in the three different GPU types that we have at our disposal, 1080Ti, Titan Xp and, Titan X.

6.5. Qualitative Results

Figure 4 shows examples of the predicted masks using our approach, FaSTGAN. The first column displays the reference mask M_n and the rest of the columns display the segmented mask by our method in the following frames.

Note that even when the input frame is corrupted by occlusions, changes of appearance and dynamic background, our method remains robust.

7. Conclusions

We have presented FaSTGAN, which, to the best of our knowledge, is the first real-time approach for semi-supervised video object segmentation running at 32 fps and yielding high-quality segmentation masks. FaSTGAN’s accuracy is on a par with previous state of the art optimized for speed but it runs at a much higher frame rate.

To achieve this, we have designed a novel GAN architecture made of a relatively small regressor and two critics that enforce spatio-temporal consistency over finite temporal windows during training. At test time, the critics are removed, leading to a simple but robust regressor that does not require fine-tuning nor post-processing operations when applied to new sequences with unseen objects. This opens a wide range of potential applications in the near future for real-time video analysis and video editing.

Acknowledgments

The authors would like to thank Fabian Mentzer and Kevis-Kokitsi Maninis for the insightful discussions and proofreading of this manuscript. This work is supported in part by the Swiss Commission for Technology and Innovation (CTI, Grant No. 19015.1 PFES-ES, NeGeVA), by the Spanish Ministry of Science and Innovation under projects HuMoUR TIN2017-90086-R, ColRobTransp DPI2016-78957 and María de Maeztu Seal of Excellence MDM-2016-0656; and by the EU project AEROARMS ICT-2014-1-644271. We also thank NVidia for hardware donation under the GPU Grant Program.

References

- [1] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein gan. *ICML*, 2017.
- [2] L. Bao, B. Wu, and W. Liu. Cnn in mrf: Video object segmentation via inference in a cnn-based higher-order spatio-temporal mrf. In *CVPR*, 2018.
- [3] A. Brock, J. Donahue, and K. Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [4] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. One-shot video object segmentation. In *CVPR*, 2017.
- [5] S. Chandra, C. Couprie, and I. Kokkinos. Deep spatio-temporal random fields for efficient video segmentation. In *CVPR*, 2018.
- [6] J. Chang, D. Wei, and J. W. Fisher III. A video representation using temporal superpixels. In *CVPR*, 2013.
- [7] Y. Chen, J. Pont-Tuset, A. Montes, and L. Van Gool. Blazingly fast video object segmentation with pixel-wise metric learning. In *CVPR*, 2018.
- [8] J. Cheng, Y.-H. Tsai, W.-C. Hung, S. Wang, and M.-H. Yang. Fast and accurate online video object segmentation via tracking parts. In *CVPR*, 2018.
- [9] J. Cheng, Y.-H. Tsai, S. Wang, and M.-H. Yang. Segflow: Joint learning for video object segmentation and optical flow. In *ICCV*, 2017.
- [10] M.-M. Cheng, N. J. Mitra, X. Huang, P. H. Torr, and S.-M. Hu. Global contrast based salient region detection. *TPAMI*, 2015.
- [11] H. Ci, C. Wang, and Y. Wang. Video object segmentation by learning location-sensitive embeddings. In *ECCV*, 2018.
- [12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [13] H. Dong, A. Supratak, L. Mai, F. Liu, A. Oehmichen, S. Yu, and Y. Guo. TensorLayer: A versatile library for efficient deep learning development. *ACM Multimedia*, 2017.
- [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 *ICCV*, 2015.
- [15] M. Everingham, L. V. Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [17] M. Grundmann, V. Kwatra, M. Han, and I. A. Essa. Efficient hierarchical graph-based video segmentation. In *CVPR*, 2010.
- [18] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved training of wasserstein GANs. In *NIPS*, 2017.
- [19] B. Hariharan, L. B. P. Arbeláez, S. Maji, and J. Malik. Semantic contours from inverse detectors. *ICCV*, 2011.
- [20] K. He, G. Gkioxari, P. Dollar, and R. Girshick. Mask R-CNN. In *ICCV*, 2017.
- [21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [22] P. Hu, G. Wang, X. Kong, J. Kuen, and Y.-P. Tan. Motion-guided cascaded refinement network for video object segmentation. In *CVPR*, 2018.
- [23] Y.-T. Hu, J.-B. Huang, and A. Schwing. Maskrnn: Instance level video object segmentation. In *NIPS*, 2017.
- [24] Y.-T. Hu, J.-B. Huang, and A. G. Schwing. Videomatch: Matching based video object segmentation. In *ECCV*, 2018.
- [25] W.-C. Hung, Y.-H. Tsai, Y.-T. Liou, Y.-Y. Lin, and M.-H. Yang. Adversarial learning for semi-supervised semantic segmentation. In *BMVC*, 2018.
- [26] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, 2017.
- [27] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017.
- [28] S. Jain, B. Xiong, and K. Grauman. Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. *CVPR*, 2017.
- [29] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *ICLR*, 2018.
- [30] D. Kingma and J. Ba. ADAM: A method for stochastic optimization. In *ICLR*, 2015.
- [31] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *CVPR*, 2017.
- [32] X. Li and C. Change Loy. Video object segmentation with joint re-identification and attention-aware mask propagation. In *ECCV*, 2018.
- [33] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and C. Zitnick. Microsoft COCO: Common objects in context. *ECCV*, 2014.
- [34] P. Luc, C. Couprie, S. Chintala, and J. Verbeek. Semantic segmentation using adversarial networks. In *NIPS Workshop on Adversarial Training*, 2016.
- [35] J. Luiten, P. Voigtlaender, and B. Leibe. Premvos: Proposal-generation, refinement and merging for video object segmentation. In *ACCV*, 2018.

- [36] Y. Luo, Z. Zheng, L. Zheng, T. Guan, J. Yu, and Y. Yang. Macro-micro adversarial network for human parsing. In *ECCV*, 2018.
- [37] N. Maerki, F. Perazzi, O. Wang, and A. Sorkine-Hornung. Bilateral space video segmentation. In *CVPR*, 2016.
- [38] K.-K. Maninis, S. Caelles, Y. Chen, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. Video object segmentation without temporal information. *TPAMI*, 2018.
- [39] N. Mayer, E. Ilg, P. Häusser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, 2016.
- [40] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun. Large kernel matters improve semantic segmentation by global convolutional network. In *CVPR*, 2017.
- [41] F. Perazzi, A. Khoreva, R. Benenson, B. Schiele, and A. Sorkine-Hornung. Learning video object segmentation from static images. In *CVPR*, 2017.
- [42] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016.
- [43] F. Perazzi, O. Wang, M. Gross, and A. Sorkine-Hornung. Fully connected object proposals for video segmentation. In *ICCV*, 2015.
- [44] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-Hornung, and L. Van Gool. The 2017 DAVIS challenge on video object segmentation. *arXiv:1704.00675*, 2017.
- [45] A. Pumarola, A. Agudo, A. Martinez, A. Sanfeliu, and F. Moreno-Noguer. GANimation: Anatomically-aware facial animation from a single image. In *ECCV*, 2018.
- [46] A. Pumarola, A. Agudo, A. Sanfeliu, and F. Moreno-Noguer. Unsupervised person image synthesis in arbitrary poses. In *CVPR*, 2018.
- [47] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *ICLR*, 2015.
- [48] M. Saito, E. Matsumoto, and S. Saito. Temporal generative adversarial nets with singular value clipping. In *ICCV*, 2017.
- [49] J. Shi, Q. Yan, L. Xu, and J. Jia. Hierarchical image saliency detection on extended cssd. *TPAMI*, 2016.
- [50] J. Shin Yoon, F. Rameau, J. Kim, S. Lee, S. Shin, and I. So Kweon. Pixel-level matching for video object segmentation using convolutional neural networks. In *ICCV*, 2017.
- [51] N. Souly, C. Spampinato, and M. Shah. Semi and weakly supervised semantic segmentation using generative adversarial network. In *ICCV*, 2017.
- [52] D. Tsai, M. Flagg, A. Nakazawa, and J. M. Rehg. Motion coherent tracking using multi-label MRF optimization. *IJCV*, 2012.
- [53] Y.-H. Tsai, W.-C. Hung, S. Schulter, K. Sohn, M.-H. Yang, and M. Chandraker. Learning to adapt structured output space for semantic segmentation. In *CVPR*, 2018.
- [54] S. Tulyakov, M.-Y. Liu, X. Yang, and J. Kautz. MocoGAN: Decomposing motion and content for video generation. *CVPR*, 2018.
- [55] C. Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [56] P. Voigtlaender and B. Leibe. Online adaptation of convolutional neural networks for video object segmentation. In *BMVC*, 2017.
- [57] C. Vondrick, H. Pirsiavash, and A. Torralba. Generating videos with scene dynamics. In *NIPS*, 2016.
- [58] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro. Video-to-video synthesis. In *NeurIPS*, 2018.
- [59] X. Wang, A. Shrivastava, and A. Gupta. A-Fast-RCNN: Hard positive generation via adversary for object detection. In *CVPR*, 2017.
- [60] S. Wug Oh, J.-Y. Lee, K. Sunkavalli, and S. Joo Kim. Fast video object segmentation by reference-guided mask propagation. In *CVPR*, 2018.
- [61] H. Xiao, J. Feng, G. Lin, Y. Liu, and M. Zhang. Monet: Deep motion exploitation for video object segmentation. In *CVPR*, 2018.
- [62] S. Xie and Z. Tu. Holistically-nested edge detection. In *ICCV*, 2015.
- [63] N. Xu, L. Yang, Y. Fan, J. Yang, D. Yue, Y. Liang, B. Price, S. Cohen, and T. Huang. YouTube-VOS: Sequence-to-sequence video object segmentation. In *ECCV*, 2018.
- [64] N. Xu, L. Yang, Y. Fan, D. Yue, Y. Liang, J. Yang, and T. Huang. YouTube-VOS: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018.
- [65] L. Yang, Y. Wang, X. Xiong, J. Yang, and A. K. Katsaggelos. Efficient video object segmentation via network modulation. In *CVPR*, 2018.