# Interactive Video Object Segmentation in the Wild

Arnaud Bénard[*]
gifs.com
arnaud@gifs.com

Michael Gygli[*]
gifs.com
michael@gifs.com

## Abstract

*In this paper we present our system for human-in-the-loop video object segmentation. The backbone of our system is a method for one-shot video object segmentation [3]. While fast, this method requires an accurate pixel-level segmentation of one (or several) frames as input. As manually annotating such a segmentation is impractical, we propose a deep interactive image segmentation method, that can accurately segment objects with only a handful of clicks. On the GrabCut dataset, our method obtains 90% IOU with just 3.8 clicks on average, setting the new state of the art. Furthermore, as our method iteratively refines an initial segmentation, it can effectively correct frames where the video object segmentation fails, thus allowing users to quickly obtain high quality results even on challenging sequences. Finally, we investigate usage patterns and give insights in how many steps users take to annotate frames, what kind of corrections they provide, etc., thus giving important insights for further improving interactive video segmentation.*

## 1. Introduction

Interactive Object segmentation has been recently become popular to quickly edit photos. Snapchat Backdrop [2], for example, allows to change the background of portrait pictures, after outlying the foreground. In their portrait mode product, Google [1] uses object segmentation to blur out the background of a photo. The practical usage of video object segmentation, on the other hand, has previously been limited, due to its challenges. Videos are often harder to segment due to motion blur, bad composition, occlusion, *etc*. Fully automatic methods typically fail to accurately segment more challenging sequences. The other extreme, requiring user input for each frame, is impractical due to its time requirements. Thus, most research in video object segmentation adapts a semi-supervised approach. There, a subset of one or several frames in a sequence are manually segmented and used to infer the seg-



Figure 1: Overview of our method. The first frame is segmented with a handful of clicks. Then, the video snippet is segmented using OSVOS [3]. Finally, the user can optionally refine poorly segmented frames with 1-2 clicks.

mentation masks of all frames in the sequence [3, 9, 15].

All aforementioned methods assume that the training frames are manually annotated with pixel-accurate segmentations. This prevents their practical usage, as annotating a single frame typically takes more than a minute [12]. In this paper, we build upon the semi-supervised approach of [3]. But our work makes it practical by proposing a fast interactive object segmentation algorithm, which allows to segment the first frame in a few seconds (see Figure 1). In our experiments, we show that using these masks, rather than costly pixel-accurate segmentations, leads to a minimal performance degradation ($-3.2\%$ IOU).

Our method for interactive segmentation is based on [18], which uses a deep convolutional neural network that jointly uses user clicks and RGB information to segment images. The key idea of our approach is to use the current segmentation mask as an additional input. Thus, our method uses the user clicks to *refine* an initial segmentation. As our experiments show, our model offers high image segmentation quality and – more importantly – allows to
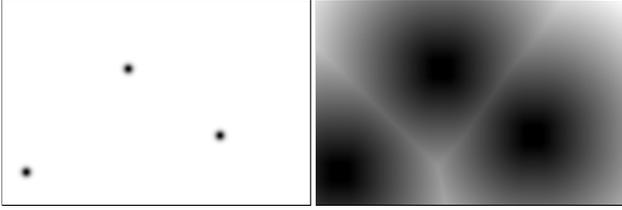
---

[*]Authors contributed equally

Figure 2: Gaussians placed at click positions *vs.* Euclidean distance map, in this case for positive clicks.

quickly correct segmentations obtained by another method, OSVOS [3] in our case. Our methodological improvements make our system usable practical and allowed us to publicly release it. By monitoring usage, we were able to obtain insights into its usage patters, something that is important for designing future video segmentation systems. We find, for example, that there are three common types of annotation and that the quality is much lower *in the wild* compared to the controlled settings that prevail in video object segmentation [14].

To summarize, this work makes the following contributions:

- A system for interactive video object segmentation based on [3].

- A novel method for deep interactive object segmentation that iteratively refines an initial segmentation by using user input in the form of clicks.

- A thorough experimental evaluation of our method and a comparison to state-of-the-art interactive segmentation methods. Our method sets a new state of the art on the GrabCut dataset [16]. We further show that it works well for providing a fast initialization for OS-VOS.

- An analysis of annotation patterns on our production data.

## 2. Method

Our approach for video object segmentation consists of two main steps: The semi-automatic annotation of the first frame and the automatic propagation of this segmentation mask to the remaining frames of the video. We will now describe these steps in turn. Finally, as video segmentation might give unsatisfactory results on challenging sequences with small objects or a lot of variation, we show how our method trivially generalizes to the case of refining an initial video segmentation result.

### 2.1. Interactive image segmentation

Interactive image segmentation is always *iterative* image segmentation, where a user provides input to *refine* the



Figure 3: Sampling corrections. Left: Current prediction and the sampled correction clicks. Right: Ground Truth mask. Positive/Negative points are sampled from the error region of the initial prediction.

current result. We propose an approach that can implicitly model that causality by providing the current segmentation mask as an input to the model, together with the user input. Specifically, we propose a deep convolutional neural network, which predicts a pixel-level segmentation as a function of the current segmentation, the user inputs and the RGB image. The current segmentation is simply a binary image provided as as an extra channel. To encode user input in the form of clicks, we create a zero-initialized channel, from which Gaussians with a small standard deviation are added or subtracted for each positive or negative click, respectively (See Figure 3). The idea of using extra channels with user information is inspired by Xu *et al.* [18], but they use two channels, each encoding the Euclidean distance to the closest foreground or background click, respectively. We choose Gaussians instead, as we find that they lead to improved localization accuracy, which is important for obtaining high-accuracy segmentations through iterative refinement. Furthermore, they allow to encode both types of clicks in a single channel. We provide a visual comparison of the two encodings in Figure 2.

#### 2.1.1 Simulating user interactions

As it is too expensive to manually collect user interactions, we follow Xu *et al.* [18] and simulate interactions instead. We iteratively sample foreground clicks that are $d_{margin} = 3$ away from the object boundary and $d_{step} = 5$ pixels away from the previous clicks. For sampling negatives we use the same 3 strategies proposed by [18]: (i) randomly sampling in a hull of $d_{hull} = 40$ around the object (ii) sampling negatives clicks from other objects in the image and (iii) sampling negatives as in (i), but such that they are maximally distant from each other, *i.e.* surround the object of interest.

In addition to that, we generate examples that correct an initial mask. For this, we sample clicks as above, run them through the model that is trained without initial masks to ob-

tain a predicted segmentation. Given these predictions, we sample $N_{corr} \in \{0, 1, 2, 3, 4, 5, 10, 20\}$ corrections clicks. Correction clicks are points that that lie in the error region. We sample them such that they are $d_{step}$ away from each other. In Figure 3 we visualize this procedure on an example. Using this kind of training examples allows the model to learn to use an initial mask, in addition to user clicks. This is what allows our model to quickly correct OSVOS results.

### 2.1.2 Neural Network architecture and training

As our network architecture we choose a ResNet-101 [8] model that is adapted to segmentation via the introduction of atrous convolutions and a pyramid scene parsing module [4]. We use weights from a model that is pre-trained on ImageNet and fine-tuned on PASCAL for semantic segmentation. The network is trained to predict the foreground probability for each pixel, via a standard cross-entropy loss. Deep neural networks cannot encode hard constraints and thus might not assign a probability of 1 for a pixel that the user labelled as foreground (and equivalently for background labelled pixels). Thus, we further clamp these pixels to be foreground/background, as in [18].

### 2.1.3 CRF Post-processing

Dominant CNN architectures, including the one we use, contain max-pooling and down-sampling layers. Thus, predictions are made at an lower resolution (8x in our case) and thus often poorly localized. To improve localization we propose to refine the initial predictions with a fully connected Conditional Random Field (CRF) [11]:

$$E(\mathbf{x}) = \sum_i \theta_i(x_i) + \sum_{ij} \theta_{ij}(x_i, x_j), \qquad (1)$$

with unary potential $\theta_i(x_i) = -\log P(x_i)$, where $P(x_i)$ is the foreground probability of pixel $i$, as predicted by the neural network. For the pairwise potential $\theta_{ij}(x_i, x_j)$ we use two Gaussian potentials as in [11]. One, favouring nearby pixels with similar colors to take the same label and the other, favouring smoothness, *i.e.* encouraging to give neighbouring pixels the same label. See [11] for more information.

### 2.2. Video segmentation

For segmenting the full video, given the first frame, we rely on one-shot video segmentation. There, the task is to segment the full video, given the segmentation of the first frame. Specifically, we use OSVOS [3]. OSVOS uses a VGG-Network [17] that is fine-tuned for video object segmentation. This network has thus learnt about objectness and which objects or things are likely to be foreground or

| Method | Pascal | GrabCut |
|---|---|---|
| iFCN w/o CRF | 6.4 | 6.2 |
| iFCN | 6.3 | 4.0 |
| Ours w/o init. mask | 5.5 | 3.9 |
| Ours w/o CRF | 4.9 | 4.7 |
| Ours | 5.6 | 3.8 |

Table 1: Ablation study.

background. To adapt the network to a particular object in a video sequence, it is fine-tuned on the segmentation mask of the first frame, such that it can learn the object's appearance. We chose OSVOS over more recent methods such as [10], due to its speed. OSVOS can be fine-tuned for 500 epochs in about 1 minute. After fine-tuning, OSVOS takes only $\approx 100ms$ per frame.

### 2.3. Progressive Refinement

Our goal is to obtain high-quality video segmentation results, even on challenging sequences. As one-shot segmentation might fail on such sequences, we need a way for the user to efficiently correct the initial result. [3] proposes a progressive refinement, where the worst segmentation mask is manually annotated and used as an additional training example for fine-tuning OSVOS. We follow this idea of correcting the worst segmentation mask, but speed up the process by relying on our interactive segmentation method presented in Section 2.1. Importantly, as our method takes the current segmentation mask as input, it provides a fast and intuitive way to *refine* an initial mask, rather than requiring to have it annotated from zero. This allows to correct bad masks with very few clicks, as our experiments in Section 3.3 show.

## 3. Experiments

We evaluate our method on three publicly available datasets Pascal [5], GrabCut [16] and DAVIS-2016 [14]. All our models are trained on 10582 images of PASCAL, augmented with the labels of SBD [7] as is common practice.

First, we evaluate our image segmentation method in Section 3.1. Then, in Section 3.2, we show how video object segmentation (OSVOS [3]) performs when initialized from a mask obtained using above method, rather than a pixel-accurate segmentation. Finally, Section 3.3 evaluates the ability to quickly correct an initial mask.

For evaluation we use use the mean Intersection over Union (IOU) [3], which is also called Jaccard index.

(a) GrabCut dataset. Top: Our predictions; Bottom: Ground Truth     (b) PASCAL dataset. Top: Our predictions; Bottom: Ground Truth

Figure 4: Results of our deep interactive object segmentation method on the GrabCut and PASCAL dataset. Our method can segment simple images with 1-2 clicks and also performs well when segmenting challenging, partially occluded objects. The rightmost image shows a failure case: It has problems segmenting thin structures such as the legs of chairs.

## 3.1. Instance Segmentation in Images

In this Section we evaluate the performance of our interactive image segmentation method of public datasets. We compare it against state-of-the-art methods, as well as our own implementation of [18], using the same DeepLab network [4], CRF inference [11] and trained with the augmented labels of [7].

### 3.1.1   Ablation study

We perform an ablation study and evaluate the effect of (i) using Gaussians instead of distance maps, (ii) the CRF refinement step and (iii) using the current mask as an additional input channel. Results are shown in Table 1.

From the table we observe that using Gaussians significantly outperforms using distance maps. This is especially pronounced when aiming for results with a high minimum IOU of 90% and not using a CRF: iFCN w/o CRF needs $6.2$ clicks, while ours w/o CRF only needs $4.7$. This matches our visual analysis of the results, where we find that iFCN does not always take user clicks into account, *i.e.* is more likely to mislabel points that were annotated by a user. We attribute this to the linearity of the distance map, which leads to more poorly localized clicks, compared to using Gaussians (*c.f.* Figure 2). The use of a CRF increases the performance on the GrabCut dataset, but decreases it on Pascal. We believe this is due to the difficulty of the Pascal images, compared to GrabCut, thus preventing gains based on simple color statistics as used in the pairwise potentials of the CRF. While using the CRF performs worse on PASCAL, we opt for using it in our method, as it improves boundary accuracy. Furthermore, we find that most users of our tool select large foreground objects such as people or animals, rather than small and often partially occluded objects, as they are common in PASCAL. Thus, we believe the GrabCut dataset is closer to our real-world data.

Finally, we observe that providing the current segmentation mask as an additional channel leads to no significant performance difference. It however allows more quickly

| Method | Pascal@85% | GrabCut@90% |
|---|---|---|
| GrabCut [16] | 15.06 | 11.10 |
| iFCN [18] | 6.88 | 6.04 |
| RIS-Net [6] | 5.12 | 5.00 |
| DEXTR [12] | **4.0** | 4.0 |
| Ours | 5.6 | **3.8** |

Table 2: The mean number of clicks required to achieve a certain IOU on different datasets by various algorithms. The best results are emphasized in **bold**.

correct an initial segmentation result, as we show in Section 3.3, which is our main motivation for using it in our model.

### 3.1.2   Comparison to the State of the Art

We compare our method against state-of-the-art methods for interactive image segmentation in Table 2. In particular, we compare against [18], on which our method is based, and the recent extension of [6]. We also compare against DEXTR [12], which uses similar techniques, but uses extreme clicks [13] as input, rather than foreground/background clicks.

**Results.**   As can be seen from the table, our method outperforms all previous methods on the GrabCut dataset. We attribute this to the use of a powerful CNN architecture (ResNet) and the use of Gaussians to encode user input. Furthermore, while [12] always needs at least four clicks, our method only needs 1-2 clicks on simple instances. On the Pascal dataset, our method is outperformed by [6, 12]. The objects in this dataset are small and thus challenging to segment, given their low resolution. Indeed, both of the superior methods use techniques to crop the region of interest, thus allowing them to obtain a high resolution input even for extremely small objects. In Figure 4 we show qualitative results. Our method often needs only 1-2 clicks to segment
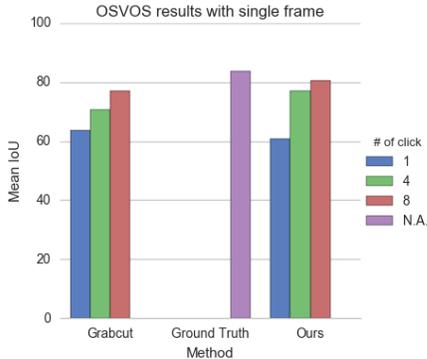
Figure 5: Comparison of OSVOS performance (mean IoU) with a single input mask generated by GrabCut, our method or taken from the ground truth.

large objects and non-occluded objects and is also able to correctly segment objects that are partially occluded.

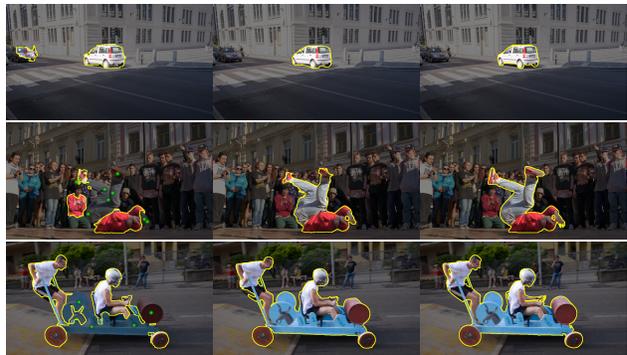## 3.2. Video object segmentation

The goal of this experiment is to measure the end-to-end accuracy of our product which outputs a full video segmentation, given a few user clicks on the first frame. To do this, we evaluate the ability of our deep image segmentation to provide an initial mask for one-shot video object segmentation.

In the experiment, we evaluate the performance of OS-VOS initialized with a first frame segmentation that is (i) the ground truth, (ii) obtained using GrabCut or (iii) obtained with our deep image segmentation. We use a different number of clicks $N_{clicks} \in \{1, 4, 8\}$ as an input to the interactive segmentation methods and evaluate the performance of OSVOS on the DAVIS-2016 [14] validation set. For GrabCut, we provide a bounding box, in addition to the user clicks.

**Results.** In Figure 5 we compare the performance obtained by using the ground truth mask against the performance when using a mask obtained with GrabCut or our method. Our method significantly outperforms using Grab-Cut, when then user provides 4 or 8 clicks. GrabCut performs better for one click, but this is because GrabCut is initialized from the ground truth bounding box, i.e. has more information. Despite this additional information provided to the GrabCut algorithm, our method performs better once the user provides more clicks. Given 8 clicks, our method obtains a mean IOU of $80.6$, while using the ground truth has an IOU of $83.8$. Thus, using a fast interactive segmentation method, rather than a tedious and time-consuming pixel-accurate segmentation, is preferable from a user experience point of view.

| Method | OSVOS | 1 click | 4 clicks | 10 clicks |
|---|---|---|---|---|
| GrabCut | 50.4% | 46.6% (-3.7) | 53.5% (+3.2) | 68.8% (+18.4) |
| iFCN | 50.4% | 55.7% (+5.3) | 71.3% (+20.9) | 79.9% (+29.5) |
| Ours | 50.4% | 63.8% (+13.4) | 75.7% (+25.4) | 82.2% (+31.8) |

Table 3: Correction of bad OSVOS masks. Our method improves the masks significantly, even for few clicks, while GrabCut and iFCN need more clicks to obtain a similar improvement.



(a) OSVOS result and corrections    (b) Our refined results    (c) Ground Truth

Figure 6: OSVOS refinement results using our interactive refinement method.

We also evaluated the performance when providing input for additional frames. We however find that adding more frames does not improve the performance significantly.

## 3.3. Segmentation Refinement

Given that an initial result obtained with OSVOS might be of unsatisfactory quality, it is important to provide a way for making intuitive and efficient corrections. In this section we thus evaluate the efficiency of our method in correcting an initial segmentation result. For this, we use the DAVIS-2016 [14] validation set. We select the worst segmentation mask per sequence, as obtained by OSVOS, and use different methods to iteratively correct the result. Thereby we compare our method, GrabCut [16] and our implementation of iFCN [18].

As GrabCut performs poorly without a bounding box (pixels that are certain background), we heuristically create a bounding box, which is created such that it includes all user foreground clicks and all probable foreground (the current segmentation mask). Without this trick, using Grab-Cut degrades the result compared to the initialization, even when using 10 correction clicks. To initialize [18], we use the center of the largest foreground blob as obtained by OS-VOS as a foreground click. We do not make it a hard con-

(a) Clicks      (b) Strokes      (c) Highlight

Figure 7: Showcase of the three types of annotations we observed from analyzing our users' behavior.

straint to allow the model to recover, should this initial click be incorrect.

**Results.** We show quantitative results in Table 3 and visual results of our method in Figure 6. As can be seen from the table, our method outperforms both baselines. The performance difference is especially prominent for few clicks: Given a single click, iFCN [18] has only a moderate improvement of 5.3% and GrabCut even decreases the performance. Our method, on the other hand, increases the IoU of the masks by as much as 13.4%. Thus, this highlights the importance of initializing a method with the existing result.

## 4. Analysis

Since we released the tool for video segmentation, termed *sticker editor*, our users have created hundreds of segmentations from videos and images. We call these segmentation *animated stickers*. In this section, we analyze a subset (437 stickers) of our production data.

### 4.1. Usage patterns

The users' annotation patterns fall into three categories (as illustrated in Figure 7): independent clicks, strokes and highlighting. Therefore, our interactive image segmentation model needs to infer accurate results with the former patterns as input. In practice, to make our model more robust to the various kind of inputs, we thus train our model with a combination of simulated clicks and strokes. Without that, we find that our model does not perform well when a user provides strokes rather than clicks.

It is not trivial to categorize user inputs into clicks, strokes or highlighting. Thus, rather than analyzing the number of annotated pixels, we use the number of refinement iterations to illustrate the amount of efforts the users invested. With a median of 4 iterative steps on both the first (Figure 8) and the refined frames (Figure 9) , we can observe that users tend to refine their masks only a few times.

After getting the first result, the user has the option to annotate more frames to refine the initial result. We find that only 15% of our users invest time to do a refinement. We can attribute this to our current user flow focused on annotating one frame.
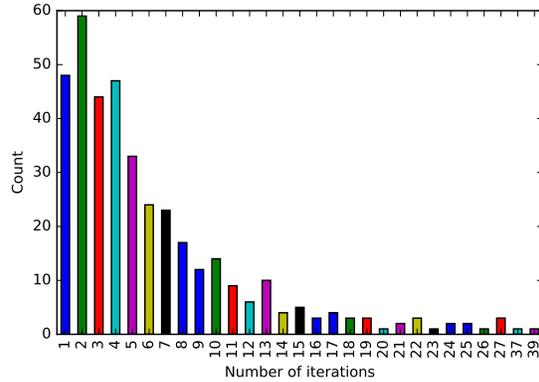


Figure 8: Histogram of the number of iterations for segmenting the first frame.
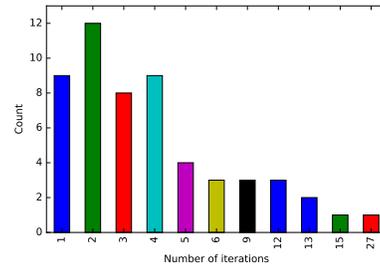


Figure 9: Histogram of the number of iterations per video frame when refining OSVOS results.
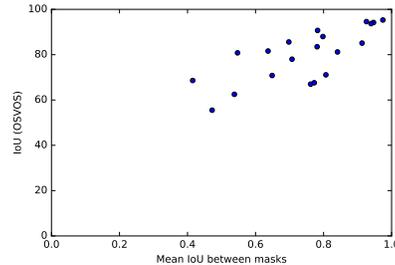


Figure 10: IOU overlap *vs*. OSVOS performance. The two are correlated, indicating that stronger frame to frame changes of the masks negatively affect OSVOS performance.

### 4.2. Future improvements

From analyzing the user data, we find that people annotate inaccurately. While we currently treat user inputs as hard constraints, these inaccuracies suggest that the user input should be used as an indication only.

We find that our users mostly create stickers of people's faces and of pets such as cats and dogs. This indicates

that pre-training OSVOS with a dataset of people and pets would improve the quality of a large portion of stickers. We also did an analysis of the performance of OSVOS on the DAVIS dataset. Thereby, we noticed a correlation between the mean IoU between frame masks and OSVOS precision (Figure 10). This shows which shows that fast moving objects are harder to segment. From a user experience point of view, it could make sense to predict the difficulty of a sequence due to motion or other factors (*e.g.* using optical flow). Such a difficulty estimate would allow to guide users towards easier sequences and thus help them make better stickers.

## 5. Conclusion

In this paper, we have presented our method for fast and interactive video object segmentation. Towards our goal of making video object segmentation practical in the wild, we have made several technical contributions. We have proposed a method to speed up the annotation of the first frame and have introduced a way for correcting mistakes of the video object segmentation, via a novel interactive segmentation model. We have empirically evaluated the different components and the full system. Our experiments showed that our model is competitive or superior to existing methods for interactive segmentation. We also showed that using masks obtained with our interactive segmentation method, rather than perfect pixel-accurate masks, affects the OSVOS results only minimally. Finally, we have provided insights into usage patterns. We have shown, for example, that most users tend to do very few iterations to annotate the first frame, thus highlighting the importance of a strong interactive segmentation model. We hope that our work and the insights we provide will spur further research on making better video object segmentation tools.

## References

[1] Portrait mode on the Pixel 2 and Pixel 2 XL smartphones. `https://research.googleblog.com/2017/10/portrait-mode-on-pixel-2-and-pixel-2-xl.html`. Accessed: 2017-12-057.

[2] Snapchat Backdrop. `https://support.snapchat.com/en-US/article/overview-backdrop`. Accessed: 2017-12-05.

[3] 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.

[4] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. DeepLab : Semantic Image Segmentation with Deep Convolutional Nets , Atrous Convolution , and Fully Connected CRFs. *PAMI*, 2017.

[5] M. Everingham, L. Van Gool, C. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge 2012 (voc2012) results (2012). 2011.

[6] J. Hao Liew, Y. Wei, W. Xiong, S.-H. Ong, and J. Feng. Regional interactive image segmentation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[7] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse detectors. In *ICCV*, 2011.

[8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[9] A. Khoreva, R. Benenson, E. Ilg, T. Brox, and B. Schiele. Lucid data dreaming for object tracking. *CoRR*, abs/1703.09554, 2017.

[10] A. Khoreva, R. Benenson, E. Ilg, T. Brox, and B. Schiele. Lucid data dreaming for object tracking. *arXiv preprint arXiv:1703.09554*, 2017.

[11] P. Krähenbühl and V. Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *NIPS*, 2011.

[12] K.-K. Maninis, S. Caelles, J. Pont-Tuset, and L. Van Gool. Deep extreme cut: From extreme points to object segmentation. *arXiv preprint arXiv:1711.09081*, 2017.

[13] D. P. Papadopoulos, J. R. Uijlings, F. Keller, and V. Ferrari. Extreme clicking for efficient object annotation. *arXiv preprint arXiv:1708.02750*, 2017.

[14] 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 *Computer Vision and Pattern Recognition*, 2016.

[15] 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.

[16] C. Rother, V. Kolmogorov, and A. Blake. Grabcut: Interactive foreground extraction using iterated graph cuts. In *ACM transactions on graphics (TOG)*. ACM, 2004.

[17] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[18] N. Xu, B. Price, S. Cohen, J. Yang, and T. S. Huang. Deep interactive object selection. In *CVPR*, 2016.