

Deep Extreme Cut: From Extreme Points to Object Segmentation

K.-K. Maninis* S. Caelles* J. Pont-Tuset L. Van Gool
 Computer Vision Lab, ETH Zürich, Switzerland

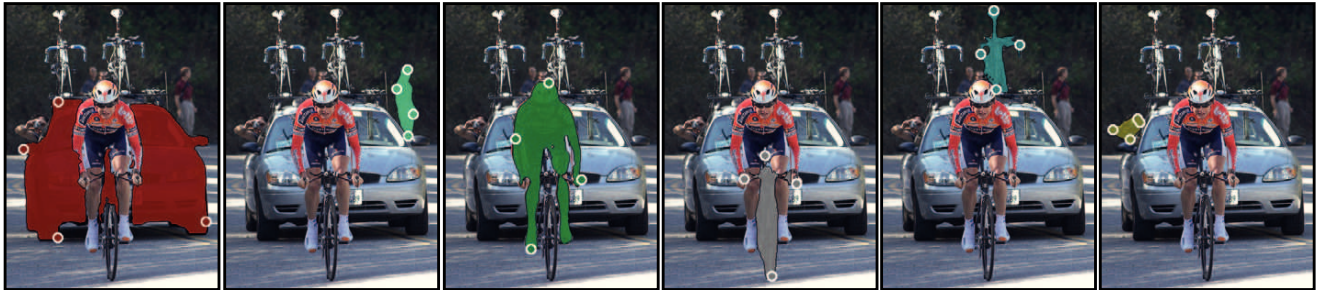


Figure 1. **Example results of DEXTR:** The user provides the extreme clicks for an object, and the CNN produces the segmented masks.

Abstract

This paper explores the use of extreme points in an object (left-most, right-most, top, bottom pixels) as input to obtain precise object segmentation for images and videos. We do so by adding an extra channel to the image in the input of a convolutional neural network (CNN), which contains a Gaussian centered in each of the extreme points. The CNN learns to transform this information into a segmentation of an object that matches those extreme points.

We demonstrate the usefulness of this approach for guided segmentation (grabcut-style), interactive segmentation, video object segmentation, and dense segmentation annotation. We show that we obtain the most precise results to date, also with less user input, in an extensive and varied selection of benchmarks and datasets. All our models and code are publicly available on <http://www.vision.ee.ethz.ch/~cvlsegmentation/dextr/>.

1. Introduction

Deep learning techniques have revolutionized the field of computer vision since their explosive appearance in the ImageNet competition [34], where the task is to classify images into predefined categories, that is, algorithms produce one label for each input image. Image and video segmentation, on the other hand, generate dense predictions where each pixel receives a (potentially different) output classification. Deep learning algorithms, especially Convolutional Neural Networks (CNNs), were adapted to this scenario by

removing the final fully connected layers to produce dense predictions.

Supervised techniques, those that train from manually-annotated results, are currently the best performing in many public benchmarks and challenges [43, 12]. In the case of image and video segmentation, the supervision is in the form of dense annotations, i.e. each pixel has to be annotated in an expensive and cumbersome process. Weakly-supervised techniques, which train from incomplete but easier-to-obtain annotations, are still significantly behind the state of the art. Semi-automatic techniques, which need a human in the loop to produce results, are another way of circumventing the expensive training annotations but need interaction at test time, which usually comes in the form of a bounding box [7] or scribbles [20] around the object of interest. How to incorporate this information at test time without introducing unacceptable lag, is also a challenge.

This paper tackles all these scenarios in a unified way and shows state-of-the-art results in all of them in a variety of benchmarks and setups. We present Deep Extreme Cut (DEXTR), that obtains an object segmentation from its four extreme points [25]: the left-most, right-most, top, and bottom pixels. Figure 1 shows an example result of our technique along with the input points provided.

In the context of semi-automatic object segmentation, we show that information from extreme clicking results in more accurate segmentations than the ones obtained from bounding-boxes (PASCAL, COCO, Grabcut) in a Grabcut-like formulation. DEXTR outperforms other methods using extreme points or object proposals (PASCAL), and provides a better input to video object segmentation (DAVIS 2016,

*First two authors contributed equally

DAVIS 2017). DEXTR can also incorporate more points beyond the extreme ones, which further refines the quality (PASCAL).

DEXTR can also be used to obtain dense annotations to train supervised techniques. We show that we obtain very accurate annotations with respect to the ground truth, but more importantly, that algorithms trained on the annotations obtained by our algorithm perform as good as when trained from the ground-truth ones. If we add the cost to obtain such annotations into the equation, then training using DEXTR is significantly more efficient than training from the ground truth for a given target quality.

We perform an extensive and comprehensive set of experiments on COCO, PASCAL, Grabcut, DAVIS 2016, and DAVIS 2017, to demonstrate the effectiveness of our approach. All code, pre-trained models and pre-computed results used in this project are publicly available on <http://www.vision.ee.ethz.ch/~cvlsegmentation/dextr/>.

2. Related Work

Weakly Supervised Signals for Segmentation: Numerous alternatives to expensive pixel-level segmentation have been proposed and used in the literature. Image-level labels [27], noisy web labels [1, 16] and scribble-level labels [20] are some of the supervisory signal that have been used to guide segmentation methods. Closer to our approach, [3] employs point-level supervision in the form of a single click to train a CNN for semantic segmentation and [26] uses central points of an imaginary bounding box to weakly supervise object detection. Also related to our approach, [7, 17] train semantic segmentation methods from box supervision. Recently, Papadopoulos et al. proposed a novel method for annotating objects by extreme clicks [25]. They show that extreme clicks provide additional information to a bounding box, which they use to enhance GrabCut-like object segmentation from bounding boxes. Different than these approaches, we use extreme clicking as a form of guidance for deep architectures, and show how this additional information can be used to further boost accuracy of segmentation networks, and help various applications.

Instance Segmentation: Several works have tackled the task of grouping pixels by object instances. Popular grouping methods provide instance segmentation in the form of automatically segmented object proposals [14, 31]. Other variants provide instance-level segmentation from a weak guiding signal in the form of a bounding box [33]. Accuracy for both groups of methods has increased by recent approaches that employ deep architectures trained on large dataset with strong supervisory signals, to learn how to produce class-agnostic masks from patches [29, 30], or from bounding boxes [41]. Our approach relates to the second

group, since we utilize information from extreme clicks to group pixels of the same instance, with higher accuracy.

Interactive Segmentation from points: Interactive segmentation methods have been proposed in order to reduce annotation time. In this context, the user is asked to gradually refine a method by providing additional labels to the data. Grabcut [33] is one of the pioneering works for the task, segmenting from bounding boxes by gradually updating an appearance model. Our method relates with interactive segmentation using points as the supervisory signal. Click Carving [15] interactively updates the result of video object segmentation by user-defined clicks. Recent methods use these ideas in the pipeline of deep architectures. iFCN [42] guides a CNN from positive and negative points acquired from the ground-truth masks. RIS-Net [10] build on iFCN to improve the result by adding local context. Our method significantly improves the results by using 4 points as the supervisory signal: the extreme points.



3. Method

3.1. Extreme points

One of the most common ways to perform weakly supervised segmentation is drawing a bounding box around the object of interest [4, 39, 33, 18]. However, in order to draw the corners of a bounding box, the user has to click points outside the object, drag the box diagonally, and adjust it several times to obtain a tight, accurate bounding box. This process is cognitively demanding, with increased error rates and labelling times [25].

Recently, Papadopoulos et al. [25] have shown a much more efficient way of obtaining a bounding box using extreme clicks, spending on average 7.2 seconds instead of 34.5 seconds required for drawing a bounding box around an object [36]. They show that extreme clicking leads to high quality bounding boxes that are on par with the ones obtained by traditional methods. These extreme points belong to the top, bottom, left-most and right-most parts of the object. Extreme-clicking annotations by definition provide more information than a bounding box; they contain four points that are on the boundary of the object, from which one can easily obtain the bounding-box. We use extreme points for object segmentation leveraging their two main outcomes: the points and their inferred bounding box.

3.2. Segmentation from Extreme Points

The overview of our method is shown in Figure 2. The annotated extreme points are given as a guiding signal to the input of the network. To this end, we create a heatmap with activations in the regions of extreme points. We center a 2D Gaussian around each of the points, in order to create a single heatmap. The heatmap is concatenated with the RGB channels of the input image, to form a 4-channel input

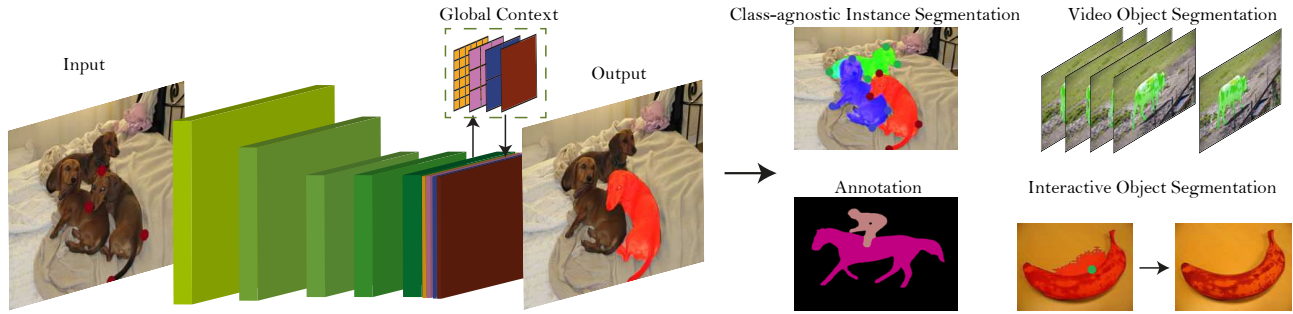


Figure 2. **Architecture of DEXTR:** Both the RGB image and the labeled extreme points are processed by the CNN to produce the segmented mask. The applicability of this method is illustrated for various tasks: Instance, Semantic, Video, and Interactive segmentation.

for the CNN. In order to focus on the object of interest, the input is cropped by the bounding box, formed from the extreme point annotations. To include context on the resulting crop, we relax the tight bounding box by several pixels. After the pre-processing step that comes exclusively from the extreme clicks, the input consists of an RGB crop including an object, plus its extreme points.

We choose *ResNet-101* [13] as the backbone of our architecture, as it has been proven successful in a variety of segmentation methods [6, 12]. We remove the fully connected layers as well as the max pooling layers in the last two stages to preserve acceptable output resolution for dense prediction, and we introduce atrous convolutions in the last two stages to maintain the same receptive field. After the last *ResNet-101* stage, we introduce a pyramid scene parsing module [43] to aggregate global context to the final feature map. Initializing the weights of the network from pre-training on ImageNet has been proven beneficial for various tasks [22, 40, 12]. For most experiments, we use the provided Deeplab-v2 model pre-trained on ImageNet, and fine-tuned on PASCAL for semantic segmentation.

The output of the CNN is a probability map representing whether a pixel belongs to the object that we want to segment or not. The CNN is trained to minimize the standard cross entropy loss, which takes into account that different classes occur with different frequency in a dataset:

$$\mathcal{L} = \sum_{j \in Y} w_{y_j} C(y_j, \hat{y}_j), \quad j \in 1, \dots, |Y| \quad (1)$$

where w_{y_j} depends on the label y_j of pixel j . In our case we define w_{y_j} with $y_j \in \{0, 1\}$ as the inverse normalized frequency of labels inside the minibatch. $C(\cdot)$ indicates the standard cross-entropy loss between the label and the prediction \hat{y}_j . The balanced loss has proven to perform very well in boundary detection [40, 23], where the majority of the samples belong to the background class. We note that our method is trained from strong mask-level supervision, on publicly available datasets, using the extreme points as a guiding signal to the network.

In order to segment an object, our method uses a object-centered crop, therefore there is a much higher number of samples belonging to the foreground than to the background

and the use of a balanced loss proves to be beneficial.

Alternatives for each of the components used in our final model have been studied in an ablation analysis, and a detailed comparison can be found in Section 4.2.

3.3. Use cases for DEXTR

Class-agnostic Instance Segmentation: One application of DEXTR is class-agnostic instance segmentation. In this task, we click on the extreme points of an object in an image, and we obtain a mask prediction for it. The selected object can be of any class, as our method is class agnostic.

In Section 4.3, we compare our method with the state of the art in two different datasets, PASCAL and Grabcut, where we improve current results. We also analyse the generalization of our method to other datasets and to unseen categories. We conclude positive results in both experiments: the performance drop for testing on a different dataset than the one used for training is very small and the result achieved is the same whether the class has been seen during training or not.

Annotation: The common annotation pipeline for segmentation can also be assisted by DEXTR. In this framework, instead of detailed polygon labels, the workload of the annotator is reduced to only providing the extreme points of an object, and DEXTR produces the desired segmentation. In this pipeline, the labelling cost is reduced by a factor of 10 (from 79 seconds needed for a mask, to 7.2 seconds needed for the extreme clicks) [25].

In Section 4.4, the quality of the produced masks are validated when used to train a semantic segmentation algorithm. We show that our method produces very accurate masks and the results trained on them are on par with those trained on the ground-truth annotations in terms of quality, with much less annotation budget.

Video Object Segmentation: DEXTR can also improve the pipeline of video object segmentation. We focus on the semi-supervised setting where methods use one or more masks as inputs to produce the segmentation of the whole video. Our aim is to replace the costly per pixel annotation masks by the masks produced by our algorithm after

the user has selected the extreme points of a certain object, and re-train strongly supervised state-of-the-art video segmentation architectures.

In Section 4.5, we provide results on two different dataset: DAVIS-2016 and DAVIS-2017. We conclude that state-of-the-art results can be achieved reducing the annotation time by a factor of 5. Moreover, for almost any specific annotation budget, better results can be obtained using a higher number of masks produced by our algorithm rather than expensive per-pixel annotated masks.

Interactive Object Segmentation: The pipeline of DEXTR can also be used in the frame of interactive segmentation from points [42, 41]. We work on the case where the user labels the extreme points of an object, but is nevertheless not satisfied with the obtained results. The natural thing to do in such case is to annotate an extra point (not extreme) in the region that segmentation fails, and expect for a refined result. **Given the nature of extreme points, we expect that the extra point also lies in the boundary of the object.**

To simulate such behaviour, we first train DEXTR on a first split of a training set of images, using the 4 extreme points as input. For the extra point, we infer on an image of the second split of the training set, and compute the accuracy of its segmentation. If the segmentation is accurate (eg. $IoU \geq 0.8$), the image is excluded from further processing. In the opposite case ($IoU < 0.8$), **we select a fifth point in the erroneous area.** To simulate human behaviour, we perturbate its location and we train the network with 5 points as input. Results presented in Section 4.6 indicate that it is possible to recover performance on the difficult examples, by using such interactive user input.

4. Experimental Validation

Our method is extensively validated on five publicly available databases: PASCAL [8], COCO [21], DAVIS 2016 [28], DAVIS 2017 [32], and Grabcut [33], for various experimental setups that show its applicability and generalization capabilities. We use DEXTR trained on PASCAL (augmented by the labels of SBD [11] following the common practice - 10582 images), unless indicated differently. Some implementation details are given in Section 4.1. We then perform an ablation study to separately validate all components of our method in Section 4.2. Class-agnostic instance segmentation experiments from extreme points are presented in Section 4.3, whereas Sections 4.4 and 4.5 are dedicated to how DEXTR contributes to segmentation annotation and video object segmentation pipelines, respectively. Section 4.6 presents our method as an interactive segmenter from points.

4.1. Implementation Details

Simulated Extreme Points: In [25], extreme points in PASCAL were obtained by crowd-sourcing. We used their collected extreme points when experimenting on the same dataset, and collected new extreme points by humans in DAVIS 2016. To experiment on COCO, on which it was not feasible to collect extreme points by human annotators, **we simulate them by taking the extreme points of the ground-truth masks jittered randomly by up to 10 pixels.**

Training and testing details: DEXTR is trained on PASCAL 2012 Segmentation for 100 epochs or on COCO 2014 training set for 10 epochs. The learning rate is set to 10^{-8} , with momentum of 0.9 and weight decay of $5 * 10^{-4}$. A mini-batch of 5 objects is used for PASCAL, whereas for COCO, due to the large size of the database, we train on 4 GPUs with an effective batch size of 20. Training on PASCAL takes approximately 20 hours on a Nvidia Titan-X GPU, and 5 days on COCO. **Testing the network is fast, requiring only 80 milliseconds.**

4.2. Ablation Study

The following sections show a number of ablation experiments in the context of class-agnostic instance segmentation to quantify the importance of each of the components of our algorithm and to justify various design choices. Table 2 summarizes these results. We use PASCAL VOC 2012 val set for the evaluation.

Architecture: We use *ResNet-101* as the backbone architecture, and compare two different alternatives. The first one is a straightforward fully convolutional architecture (Deeplab-v2 [6]) where the fully connected and the last two max pooling layers are removed, and the last two stages are substituted with dilated (or atrous) convolutions. This keeps the size of the prediction in reasonable limits ($8\times$ lower than the input). We also tested a region-based architecture, similar to Mask R-CNN [12], with a re-implementation of the ResNet-101-C4 variant [12], which uses the fifth stage (C5) for regressing a mask from the Region of Interest (RoI), together with the re-implementation of the RoI-Align layer. For more details please refer to [12]. In the first architecture, the input is a patch around the object of interest, whereas in the latter the input is the full image, and cropping is applied at the RoI-Align stage. Deeplab-v2 performs +3.9% better. We conclude that the output resolution of ResNet-101-C4 (28×28) is inadequate for the level of detail that we target.

Bounding boxes vs. extreme points: **We study the performance of Deeplab-v2 as a foreground-background classifier given a bounding box compared to the extreme points.** In the first case, the input of the network is the cropped image around the bounding box plus a margin of 50 pixels

to include some context. In the second case, the extreme points are fed together in a fourth channel of the input to guide the segmentation. Including extreme points to the input **increases performance by +3.1%**, which suggest that they are a source of very valuable information that the network uses additionally to guide its output.

Loss: For the task of class-agnostic instance segmentation, we compare two binary losses, i.e. the standard cross-entropy and a class-balanced version of it, where the loss for each class in the batch is weighted by its inverse frequency. Class-balancing the loss gives more importance to the less frequent classes, and has been successful in various tasks [40, 5]. **DEXTR also performs better when the loss is balanced, leading to a performance boost of +3.3%.**

Full image vs. crops: Having the extreme points annotated allows for focusing on specific regions in an image, cropped by the limits specified by them. In this experiment, we compare how beneficial it is to focus on the region of interest, rather than processing the entire image. To this end, we crop the region surrounded by the extreme points, relaxing it by 50 pixel for increased context and compare it against the full image case. **We notice that cropping increases performance by +7.9%, and is especially beneficial for the small objects of the database.** This could be explained by the fact that cropping eliminates the scale variation on the input. Similar findings have been reported for video object segmentation by [19].

Atrous spatial pyramid (ASPP) vs. pyramid scene parsing (PSP) module: Pyramid Scene Parsing Network [43] steps on the Deeplab-v2 [6] architecture to further improve results on semantic segmentation. Their main contribution was a global context module (PSP) that employs global features together with the local features for dense prediction. We compare the two network heads, the original ASPP [6], and the recent PSP module [43] for our task. The increased results of the PSP module (+2.3%) indicate that the PSP module builds a global context that is also useful in our case.

Manual vs. simulated extreme points: In this section we analyze the differences between the results obtained by DEXTR when we input either human-provided extreme points or our simulated ones, to check that the conclusions we draw from the simulations will still be valid in a realistic use case with human annotators. We do so in the two datasets where we have *real* extreme points from humans. The first one is a certain subset of PASCAL 2012 Segmentation and SBD (5623 objects) with extreme points from [25], which we refer to as PASCAL_{EXT} and DAVIS 2016, for which we crowdsourced the extreme point annotations. The annotation time for the latter (average of all 1376 frames

Method	PASCAL _{EXT}	DAVIS 2016
Manual extreme points	80.1	80.9
Simulated extreme points	85.1	79.5

Table 1. **Manual vs. simulated extreme points:** Intersection over Union (IoU) of the DEXTR results when using manual or simulated extreme points as input.

of the validation set) was 7.5 seconds per frame, in line with [25] (7.2 s. per image). Table 1 shows that the results are indeed comparable when using both type of inputs. The remainder of the paper uses the simulated extreme points except when otherwise specified.

Distance-map vs. fixed points: Recent works [42, 41, 10] that focus on segmentation from (not-extreme) points use the distance transform of positive and negative annotations as an input to the network, in order to guide the segmentation. We compare with their approach by substituting the fixed Gaussians to the distance transform of the extreme points. We notice a performance drop of -1.3%, suggesting that using fixed Gaussians centered on the points is a better representation when coupled with extreme points. In Section 4.3 we compare to such approaches, showing that extreme points provide a much richer guidance than arbitrary points on the foreground and the background of an object.

Component #1	Component #2	Gain in IoU
Region-based	Deeplab-v2	+3.9%
Bounding Boxes	Extreme Points	+3.1%
Cross Entropy	Balanced BCE	+3.3%
Full Image	Crop on Object	+7.9%
ASPP	PSP	+2.3%
Fixed Points	Distance Map	-1.3%

Table 2. **Ablation study:** Comparative evaluation between different choices in various components of our system. Mean IoU over all objects in PASCAL VOC 2012 val set.

Summary Table 3 illustrates the building blocks that lead to the best performing variant for our method. All in all, we start by a Deeplab-v2 base model working on bounding boxes. We add the PSP module (+2.3%), the extreme points

Variant	IoU (%)	Gain
Full Image (Deeplab-v2 + PSP + Extreme Points)	82.6	
Crop on Object (Deeplab-v2)	85.1	+2.5%
+ PSP	87.4	+2.3%
+ Extreme Points	90.5	+3.1%
+ SBD data (Ours)	91.5	+1.0%

Table 3. **Ablation study:** Building performance in PASCAL VOC 2012 val set.

in the input of the network (+3.1%), and more annotated data from SBD (+1%) to reach maximum accuracy. The improvement comes mostly because of the guidance from extreme points, which highlights their importance for the task.

4.3. Class-agnostic Instance Segmentation

Comparison to the State of the Art in PASCAL: We compare our method against state-of-the-art class-agnostic instance segmentation methods in Table 4. DEXTR gets a boost of +6.5% with respect to using the *grabcut-based* method of [25] from extreme points.

Method	IoU
Sharpmask [30] from bounding box	69.3%
Sharpmask [30] upper bound	78.0%
[25] from extreme points	73.6%
Ours from extreme points	80.1%

Table 4. **Comparison on PASCAL_{EXT}:** IoU of our results against class-agnostic instance segmentation methods, on the objects annotated by [25] to be able to compare to them.

We then compare to two other *baselines* using SharpMask [30], the state-of-the-art object proposal technique. In the first row, we evaluate the proposal (out of 1000) whose bounding box best overlaps with the ground-truth bounding box, mimicking a naive algorithm to segment boxes from proposals. The second row shows the upper bound of SharpMask, that is, the best proposal against the ground truth, selected by an oracle. Both approaches are well below our result (-10.8% and -2.1%). Figure 3 illustrates some results obtained by our method on PASCAL.

Comparison to the State of the Art on the Grabcut dataset: We use our best PASCAL model and we test it in the Grabcut dataset [33]. This dataset contains 50 images, each with one annotated object from various categories, some of them not belonging to any of the PASCAL ones (banana, scissors, kangaroo, etc.). The evaluation metric is the error rate: the percentage of misclassified pixels within the bounding boxes provided by [18]. Table 5 shows the results, where DEXTR achieves 2.3% error rate, 1.1% below the runner up (or a 32% relative improvement).

Method	Error Rate (%)
GrabCut [33]	8.1
KernelCut [37]	7.1
OneCut [38]	6.7
[25] from extreme points	5.5
BoxPrior [18]	3.7
MILCut [39]	3.6
DeepGC [41]	3.4
Ours from extreme points	2.3

Table 5. **Comparison on the Grabcut dataset:** Error rates compared to the state-of-the-art techniques.

Generalization to unseen categories and across datasets:

Table 6 shows our results when trained on a certain dataset (first column), and tested in another one or certain categories (second column). In order to make a fair comparison, all the models are pre-trained only on Imageet [34] for image labeling and trained on the specified dataset for category-agnostic instance segmentation. The first two rows show that our technique is indeed class agnostic, since the model trained on PASCAL achieves roughly the same performance in COCO mini-val (MVal) regardless of the categories tested. The remaining rows shows that DEXTR also generalizes very well across datasets, since differences are around only 2% of performance drop.

	Train	Test	IoU
Unseen categories	PASCAL	COCO MVal w/o PASCAL classes	80.3%
	PASCAL	COCO MVal only PASCAL classes	79.9%
Dataset generalization	PASCAL	COCO MVal	80.1%
	COCO	COCO MVal	82.1%
	COCO	PASCAL	87.8%
	PASCAL	PASCAL	89.8%

Table 6. **Generalization to unseen classes and across datasets:** Intersection over union results of training in one setup and testing on another one. MVal stands for mini-val.

Generalization to background (stuff) categories: In order to verify the performance of DEXTR in “background” classes, we trained a model using the background labels of PASCAL Context [24] (road, sky, sidewalk, building, wall, fence, grass, ground, water, floor, ceiling, mountain, and tree). Qualitative results (Figure 3 last row) suggest that our method generalizes to background classes as well. Quantitatively, we achieve a mIoU of 81.75% in PASCAL-Context validation set, for the aforementioned classes.

4.4. Annotation

As seen in the previous section, DEXTR is able to generate high-quality class-agnostic masks given only extreme points as input. The resulting masks can in turn be used to train other deep architectures for other tasks or datasets, that is, we use extreme points as a way to annotate a new dataset with object segmentations. In this experiment we compare the results of a semantic segmentation algorithm trained on either the ground-truth masks or those generated by DEXTR (we combine all per-instance segmentations into a per-pixel semantic classification result).

Specifically, we train DEXTR on COCO and use it to generate the object masks of PASCAL train set, on which we train Deeplab-v2 [6], and the PSP [43] head as the semantic segmentation network. To keep training time manageable, we do not use multi-scale training/testing. We evaluate the results on the PASCAL 2012 Segmentation val

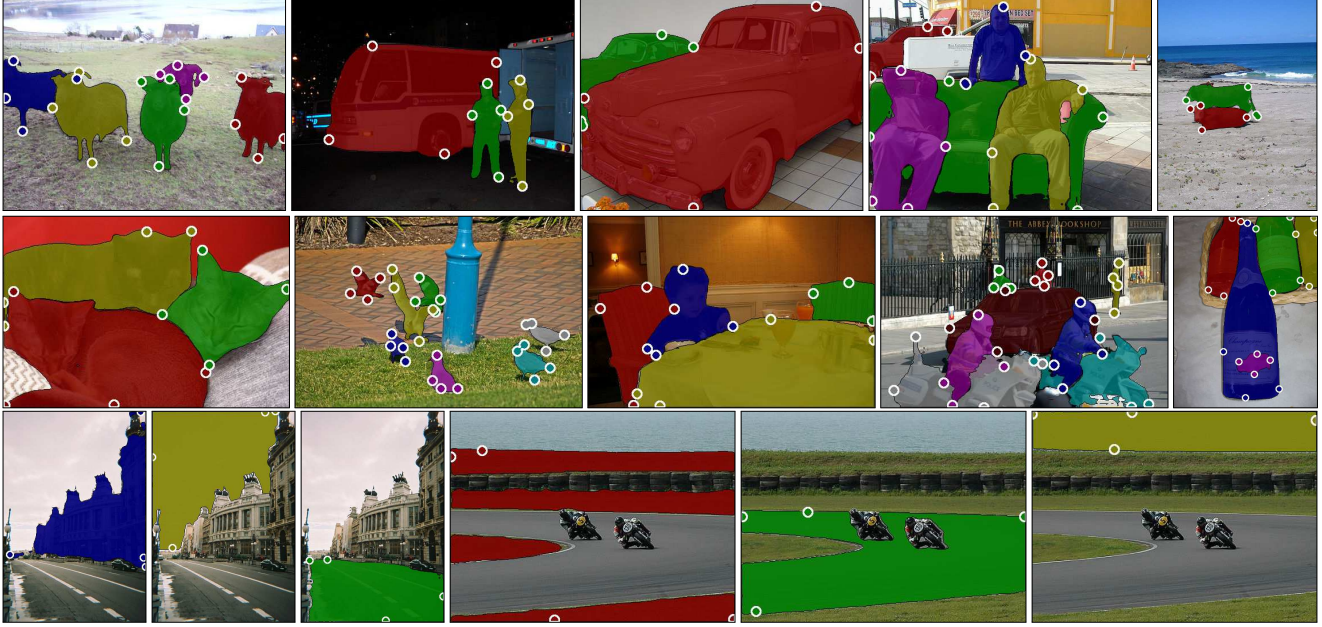


Figure 3. **Qualitative results by DEXTR on PASCAL:** Each instance with the simulated extreme points used as input and the resulting mask overlaid. The bottom row shows results on PASCAL Context stuff categories.

set, and measure performance by the standard mIoU measure (IoU per-category and averaged over categories).

Figure 4 shows the results with respect to the annotation budget (left) and the number of images (right). For completeness, we also report the results of PSPNet [43] (●) by evaluating the model provided by the authors (pre-trained on COCO, with multi-scale training and testing). The results trained on DEXTR’s masks are significantly better than those trained from the ground truth on the same budget (e.g. 70% IoU at 7-minute annotation time vs. 46% with the same budget, or 1h10 instead of 7 minutes to reach the same 70% accuracy). DEXTR’s annotations reach practically the same performance than ground truth when given the same number of annotated images.

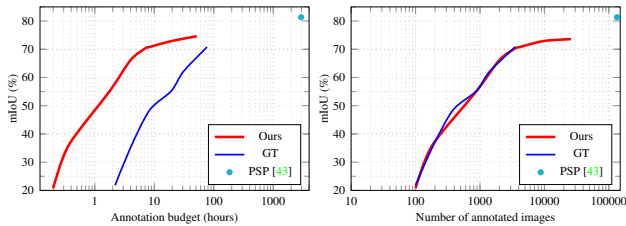


Figure 4. **Quality vs. annotation budget:** mIoU for semantic segmentation on PASCAL val set trained on our masks or the input, as a function of annotation budget (left) and the number of annotated images (right).

4.5. Video Object Segmentation

We test DEXTR also for Video Object Segmentation on the DAVIS datasets [28, 32]. We focus on the semi-supervised setting *i.e.* the mask in one or more frames of

the object that we want to segment is given as input to the algorithm, and as before we will compare the results obtained from the masks obtained by DEXTR or the ground truth having a certain annotation budget. We assume that the annotation time of the DAVIS masks is the same than that of COCO [21] (79 seconds per instance), despite the former are significantly more accurate.

We use OSVOS [5], as a state-of-the-art semi-supervised video object segmentation technique, which heavily relies on the appearance of the annotated frame, and their code is publicly available. Figure 5 (left) shows the performance of OSVOS in DAVIS 2016 [28] trained on the ground truth mask (—) or the masks generated by DEXTR from extreme points (—). We reach the same performance as using one ground-truth annotated mask with an annotation budget 5 times smaller. Once we train with more than one ground-truth annotated mask, however, even though we can generate roughly ten times more masks, we cannot achieve the same accuracy. We believe this is so because DAVIS 2016 sequences have more than one semantic instance per mask while we only annotate a global set of extreme points, which confuses DEXTR.

To corroborate this intuition, we perform the same experiment in DAVIS 2017 [32], where almost every mask contains only one instance. Figure 5 (right) shows that the performance gap with respect to using the full ground-truth mask is much smaller than in DAVIS 2016. Overall, we conclude that DEXTR is also very efficient to reduce annotation time in video object segmentation.

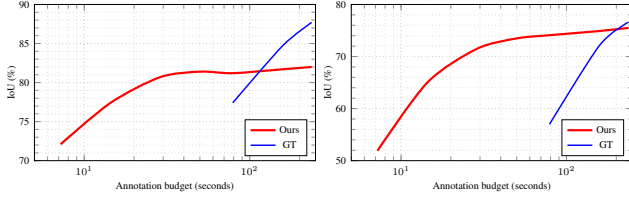


Figure 5. **Quality vs. annotation budget in video object segmentation:** OSVOS’ performance when trained from the masks of DEXTR or the ground truth, on DAVIS 2016 (left) and on DAVIS 2017 (right).

4.6. Interactive Object Segmentation

DEXTR for Interactive Segmentation: We experiment on PASCAL VOC 2012 segmentation for interactive object segmentation. We split the training dataset into two equal splits. Initially, we train DEXTR on the first split and test on the second. We then focus on the objects with inaccurate segmentations, i.e. $\text{IoU} < 0.8$, to simulate the ones on which a human - unsatisfied with the result - would mark a fifth point. The extra point would lie on the boundary of the erroneous area (false positive or false negative), which we simulate as the boundary point closest to the highest error. From the perspective of network training, this can be interpreted as Online Hard Example Mining (OHEM) [35], where one only needs to back-propagate gradients for the training examples that lead to the highest losses. Results are presented in Table 7.

Trained on	4 points	4 points-all	5 points	5 points + OHEM
IoU	59.6%	69.0%	69.2%	73.2%

Table 7. **Interactive Object Segmentation Evaluation:** Average IoU on difficult cases of PASCAL VOC 2012 validation dataset.

We first select the objects that lead to poor performance ($\text{IoU} < 0.8$) when applying the network trained on the first split. We report the average IoU on them (338 objects - 59.6%). Using the network trained further on the hard examples, with a fifth boundary point, performance increases to 73.2% (“5 points + OHEM”).

Since the increased performance is partially due to the increased amount of training data (first split + hard examples of the second split), we need to disentangle the two sources of performance gain. To this end, we train DEXTR on 4 points, by appending the hard examples of the second split to the first split of our training set (“4 points-all”).

Results suggest that DEXTR learns to handle more input information given interactively in the form of boundary clicks, to improve results of poorly segmented difficult examples (+4.2%). Interestingly, OHEM is a crucial component for improving performance: without it the network does not focus on the difficult examples (only 11% of objects of the second training split are hard examples), and fails to improve on the erroneous region indicated by the

fifth boundary point (“5 points”).

Comparison to the State of the Art: We compare against the state-of-the-art in interactive segmentation by considering extreme points as 4 clicks. Table 8 shows the number of clicks that each method needs to reach a certain performance, as well as their performance when the input is 4 clicks, in PASCAL and the Grabcut dataset. DEXTR reaches about 10% higher performance at 4 clicks than the best competing method, and reaches 85% or 90% quality with fewer clicks. This further demonstrates the enhanced performance of the CNN, when guided by extreme points.

Method	Number of Clicks		IoU (%) @ 4 clicks	
	PASCAL@85%	Grabcut@90%	PASCAL	Grabcut
GraphCut [4]	> 20	> 20	41.1	59.3
Geodesic matting [2]	> 20	> 20	45.9	55.6
Random walker [9]	16.1	15	55.1	56.9
iFCN [42]	8.7	7.5	75.2	84.0
RIS-Net [10]	5.7	6.0	80.7	85.0
Ours	4.0	4.0	91.5	94.4

Table 8. **PASCAL and Grabcut Dataset evaluation:** Comparison to interactive segmentation methods in terms of number of clicks to reach a certain quality and in terms of quality at 4 clicks.

To the meticulous reader, please note that the difference in performance of DEXTR between Table 8 (91.5%) and Table 1 (85.1%) comes from the fact that the former is on PASCAL VOC 2012 segmentation validation, so DEXTR is trained on SBD + PASCAL train, whereas the latter is on a subset of PASCAL that overlaps with train (PASCAL_{EXT}), so DEXTR is only trained on COCO.

5. Conclusions

We have presented DEXTR, a CNN architecture for semi-automatic segmentation that turns extreme clicking annotations into accurate object masks; by having the four extreme locations represented as a *heatmap* extra input channel to the network. The applicability of our method is illustrated in a series of experiments regarding semantic, instance, video, and interactive segmentation in five different datasets; obtaining state-of-the-art results in all scenarios. DEXTR can also be used as an accurate and efficient mask annotation tool, reducing labeling costs by a factor of 10.

Acknowledgments: Research funded by the EU Framework Programme for Research and Innovation Horizon 2020 (Grant No. 645331, EurEyeCase), and the Swiss Commission for Technology and Innovation (CTI, Grant No. 19015.1 PFES-ES, NeGeVA). We gratefully acknowledge support by armasuisse and thank NVidia Corporation for donating the GPUs used in this project. We thank Dim P. Papadopoulos for sharing the resources of [25] with us and Vittorio Ferrari for constructive comments and discussions.

References

- [1] E. Ahmed, S. Cohen, and B. Price. Semantic object selection. In *CVPR*, 2014. [2](#)
- [2] X. Bai and G. Sapiro. Geodesic matting: A framework for fast interactive image and video segmentation and matting. *IJCV*, 82(2):113–132, 2009. [8](#)
- [3] A. Bearman, O. Russakovsky, V. Ferrari, and L. Fei-Fei. What’s the Point: Semantic Segmentation with Point Supervision. *ECCV*, 2016. [2](#)
- [4] Y. Y. Boykov and M.-P. Jolly. Interactive graph cuts for optimal boundary & region segmentation of objects in nd images. In *ICCV*, 2001. [2](#), [8](#)
- [5] 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](#), [7](#)
- [6] 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. *T-PAMI*, 2017. [3](#), [4](#), [5](#), [6](#)
- [7] J. Dai, K. He, and J. Sun. Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *ICCV*, 2015. [1](#), [2](#)
- [8] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. [4](#)
- [9] L. Grady. Random walks for image segmentation. *T-PAMI*, 28(11):1768–1783, 2006. [8](#)
- [10] J. Hao Liew, Y. Wei, W. Xiong, S.-H. Ong, and J. Feng. Regional interactive image segmentation networks. In *ICCV*, 2017. [2](#), [5](#), [8](#)
- [11] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse detectors. In *ICCV*, 2011. [4](#)
- [12] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *ICCV*, 2017. [1](#), [3](#), [4](#)
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. [3](#)
- [14] J. Hosang, R. Benenson, P. Dollár, and B. Schiele. What makes for effective detection proposals? *T-PAMI*, 38(4):814–830, 2016. [2](#)
- [15] S. D. Jain and K. Grauman. Click carving: Segmenting objects in video with point clicks. In *HCOMP*, 2016. [2](#)
- [16] B. Jin, M. V. Ortiz Segovia, and S. Susstrunk. Webly supervised semantic segmentation. In *CVPR*, 2017. [2](#)
- [17] A. Khoreva, R. Benenson, J. Hosang, M. Hein, and B. Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *CVPR*, 2017. [2](#)
- [18] V. Lempitsky, P. Kohli, C. Rother, and T. Sharp. Image segmentation with a bounding box prior. In *ICCV*, pages 277–284, 2009. [2](#), [6](#)
- [19] X. Li, Y. Qi, Z. Wang, K. Chen, Z. Liu, J. Shi, P. Luo, C. C. Loy, and X. Tang. Video object segmentation with re-identification. *The 2017 DAVIS Challenge on Video Object Segmentation - CVPR Workshops*, 2017. [5](#)
- [20] D. Lin, J. Dai, J. Jia, K. He, and J. Sun. Scribble-sup: Scribble-supervised convolutional networks for semantic segmentation. In *CVPR*, 2016. [1](#), [2](#)
- [21] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. [4](#), [7](#)
- [22] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. [3](#)
- [23] K. Maninis, J. Pont-Tuset, P. Arbeláez, and L. V. Gool. Convolutional oriented boundaries: From image segmentation to high-level tasks. *T-PAMI*, 2017. [3](#)
- [24] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, 2014. [6](#)
- [25] D. P. Papadopoulos, J. R. Uijlings, F. Keller, and V. Ferrari. Extreme clicking for efficient object annotation. In *ICCV*, 2017. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [8](#)
- [26] D. P. Papadopoulos, J. R. Uijlings, F. Keller, and V. Ferrari. Training object class detectors with click supervision. In *CVPR*, 2017. [2](#)
- [27] D. Pathak, E. Shelhamer, J. Long, and T. Darrell. Fully convolutional multi-class multiple instance learning. In *ICLR*, 2015. [2](#)
- [28] 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. [4](#), [7](#)
- [29] P. O. Pinheiro, R. Collobert, and P. Dollár. Learning to segment object candidates. In *NIPS*, 2015. [2](#)
- [30] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollár. Learning to refine object segments. In *ECCV*, 2016. [2](#), [6](#)
- [31] J. Pont-Tuset, P. Arbeláez, J. T. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping for image segmentation and object proposal generation. *T-PAMI*, 39(1):128–140, 2017. [2](#)
- [32] 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. [4](#), [7](#)
- [33] C. Rother, V. Kolmogorov, and A. Blake. Grabcut: Interactive foreground extraction using iterated graph cuts. In *ACM TOG*, 2004. [2](#), [4](#), [6](#)
- [34] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. [1](#), [6](#)
- [35] A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, 2016. [8](#)
- [36] H. Su, J. Deng, and L. Fei-Fei. Crowdsourcing annotations for visual object detection. In *AAAI Workshops*, 2012. [2](#)
- [37] M. Tang, I. Ben Ayed, D. Marin, and Y. Boykov. Secrets of grabcut and kernel k-means. In *ICCV*, pages 1555–1563, 2015. [6](#)
- [38] M. Tang, L. Gorelick, O. Veksler, and Y. Boykov. Grabcut in one cut. In *ICCV*, pages 1769–1776, 2013. [6](#)
- [39] J. Wu, Y. Zhao, J.-Y. Zhu, S. Luo, and Z. Tu. MILcut: A sweeping line multiple instance learning paradigm for interactive image segmentation. In *CVPR*, 2014. [2](#), [6](#)

- [40] S. Xie and Z. Tu. Holistically-nested edge detection. *IJCV*, pages 1–16, 2017. 3, 5
- [41] N. Xu, B. Price, S. Cohen, J. Yang, and T. Huang. Deep grabcut for object selection. In *BMVC*, 2017. 2, 4, 5, 6
- [42] N. Xu, B. Price, S. Cohen, J. Yang, and T. S. Huang. Deep interactive object selection. In *CVPR*, 2016. 2, 4, 5, 8
- [43] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *CVPR*, 2017. 1, 3, 5, 6, 7