# Faster Training of Mask R-CNN by Focusing on Instance Boundaries☆

Roland S. Zimmermann[a,b,1], Julien N. Siems[a,c,2]

[a]*BMW Car IT GmbH, Lise-Meitner-Straße 14, 89081 Ulm, Germany*
[b]*Georg-August University of Göttingen, Friedrich-Hund-Platz 1, 37077 Göttingen, Germany*
[c]*Albert Ludwig University of Freiburg, Fahnenbergplatz, 79085 Freiburg im Breisgau, Germany*

## Abstract

We present an auxiliary task to Mask R-CNN, an instance segmentation network, which leads to faster training of the mask head. Our addition to Mask R-CNN is a new prediction head, the Edge Agreement Head, which is inspired by the way human annotators perform instance segmentation. Human annotators copy the contour of an object instance and only indirectly the occupied instance area. Hence, the edges of instance masks are particularly useful as they characterize the instance well. The Edge Agreement Head therefore encourages predicted masks to have similar image gradients to the ground-truth mask using edge detection filters. We provide a detailed survey of loss combinations and show improvements on the MS COCO Mask metrics compared to using no additional loss. Our approach marginally increases the model size and adds no additional trainable model variables. While the computational costs are increased slightly, the increment is negligible considering the high computational cost of the Mask R-CNN architecture. As the additional network head is only relevant during training, inference speed remains unchanged compared to Mask R-CNN. In a default Mask R-CNN setup, we achieve a training speed-up and a relative overall improvement of 8.1% on the MS COCO metrics compared to the baseline.

*Keywords:* Mask R-CNN, Instance Segmentation, Computer Vision, Auxiliary Task, Edge Detection Filter, Sobel Filter, Laplace Filter, Convolutional Neural Network

## 1. Introduction

Significant improvements in computer vision techniques have been made possible by the rapid progress of training Deep Convolutional Neural Networks in recent years. Application areas include image classification (Krizhevsky et al., 2012; Szegedy et al., 2015; Simonyan and Zisserman, 2015; He et al., 2016) and object detection (Girshick, 2015; Redmon et al., 2016; Liu et al., 2016). One of the most demanding computer vision tasks is instance segmentation, as it involves localizing and segmenting object instances. Recently, there have been multiple methods (Li et al., 2017; Bai and Urtasun, 2017; Liu et al., 2018; He et al., 2017) proposed to perform this task.

Another beneficial factor to the success of these Deep Learning architectures is the availability of large labeled datasets such as MS COCO (Lin et al., 2014) and the Cityscapes dataset (Cordts et al., 2016). Labeling an image dataset for instance segmentation is particularly time-consuming, because it requires segmenting all objects in a scene. It is therefore highly desirable to speed up training of an instance segmentation model to be more data efficient. In this work, we propose a conceptually straightforward addition to the Mask R-CNN (He et al., 2017) architecture which reduces training time of the mask branch.

The Mask R-CNN architecture is based on Faster R-CNN (Ren et al., 2017), which introduced an efficient Region Proposal Network (RPN) design to output bounding box proposals. The proposals are computed using a sliding window approach to make them translation invariant. A feature extractor such as ResNet (He et al., 2016), Inception (Szegedy et al., 2017) or VGGNet (Simonyan and Zisserman, 2015) is used as input to the region proposal network. The regions and features are used in the bounding box regression head, that refines the bounding box localization and the softmax classification head, which determines the instance class. This second stage is the architecture as described in Fast R-CNN (Girshick, 2015).

Mask R-CNN is a simple but effective addition to the Faster R-CNN architecture that adds a head for instance mask prediction. Using a small Fully Convolutional Neural Network (FCN) (Long et al., 2015), it can predict pixel level instance masks. Besides the mask branch, it uses a Feature Pyramid Network (FPN) backbone as proposed by Lin et al. (2017). This addition allows the network to make use of both high-resolution feature maps in the lower layers for accurate localization, as well as semantically more meaningful higher-level features, which are of lower resolution. Another contribution is ROI Align which maps arbitrarily sized spatial regions of interest in the features to a fixed spatial resolution using bilinear interpolation. This modification improves the COCO Mask metrics and enables the

---

☆Both authors contributed equally to this work and must both be cited. They are both corresponding authors.

[1]R. S. Zimmermann is with the University of Göttingen, Göttingen, Germany. This work was started when he was an intern at BMW Car IT GmbH, Ulm, Germany. Email: roland.zimmermann@stud.uni-goettingen.de

[2]J. N. Siems is with the Albert Ludwig University of Freiburg, Freiburg, Germany. This work was started when he was an intern at BMW Car IT GmbH, Ulm, Germany. Email: siemsj@cs.uni-freiburg.de

use of instance masks which require precise localization.

For the mask head, a new loss term $L_{Mask}$ has been introduced, which calculates the pixel-wise cross entropy between the predicted and target masks. The Mask R-CNN loss function

$$L_{MRCNN} = L_{Class} + L_{Box} + L_{Mask} \qquad (1)$$

is a multi-task loss based on the Faster R-CNN loss.

We propose to attach an Edge Agreement Head to the mask branch of Mask R-CNN which acts as an auxiliary task to Mask R-CNN. This head uses traditional edge detection filters such as Sobel and Laplacian kernels (Sobel and Feldman, 1973; Forsyth and Ponce, 2002) on both the predicted mask and the ground-truth mask to encourage their edges to agree. Instances in natural images are bounded by the edges that annotators use to mark the instance. Therefore, we show that encouraging the edges in the predicted and ground-truth mask to agree leads to faster training of our mask head. We argue that this is a result of the instance boundary being a robust feature to mask prediction, which can be easily propagated from the image to the mask branch.

## 2. Related Work

**Multi Task Learning.** The Edge Agreement Head acts as an auxiliary task (Ruder, 2017) to the multi-task model Mask R-CNN, which is performing both object detection and instance segmentation. Auxiliary tasks have shown to encourage models to learn robust representations of their input in a variety of applications, such as facial landmark detection (Zhang et al., 2014), natural language processing (Collobert and Weston, 2008) or steering prediction in autonomous driving (Caruana, 1997). Even seemingly unrelated tasks, e.g. weather prediction to semantic scene segmentation, can improve the model's overall performance (Liebel and Körner, 2018).
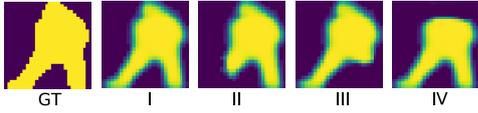
**Monocular Depth Estimation.** Godard et al. (2017) use image gradients to encourage consistency between input images and predicted disparity maps. However, the left-right disparity consistency loss does not ensure image gradients of the predicted disparity of the left and right camera to exhibit similar edge detection filter responses.

**Scene segmentation.** Chen et al. (2016) show a two-part model predicting both semantic segmentations and edges. The semantic segmentation model is based on the DeepLab model (Chen et al., 2018a) and the edge detection filter is created using intermediate convolutional filters of the DeepLab model. The task specific edge-detection on the input image is used to refine the coarse segmentation using domain transform. Our approach determines edges in fixed size, low dimensional instance mask images, for which traditional edge detection filters have been proven to be effective. Similarly, Marmanis et al. (2018) predict both semantic scene segmentation and semantic boundaries. The network responsible for predicting semantic boundaries is trained using a Euclidean loss before each pooling layer to enforce each layer to predict edges at different scales. Our approach uses predefined edge detection filters with well-known properties, which are kept constant during training, leaving us with a significantly lower additional memory footprint and computational costs.

**Edge detection.** The detection of edges has been a research topic for many decades and numerous methods have been proposed (Sobel and Feldman, 1973; Konishi et al., 2003). This field has seen large improvements due to deep learning techniques (Bertasius et al., 2015; Shen et al., 2015; Xie and Tu, 2015). Our work uses the Sobel image gradient filters proposed in (Sobel and Feldman, 1973), because it keeps the computational overhead to a minimum. Furthermore, our edge detection filters are used on $28 \times 28$ sized masks with only one channel depicting a single instance and not high-resolution color images. This significantly reduces the complexity of the problem and justifies our choice of simple edge detection filters.

**Instance segmentation.** Hayder et al. (2017) propose a model that predicts the truncated distance transform (Borgefors, 1986) of the mask, making it more resilient towards non-instance enclosing bounding box proposals. The proposed architecture for the boundary-aware instance segmentation network has many similarities to Mask R-CNN as they are both based on the Faster R-CNN architecture by Ren et al. (2017). However, they achieved lower results on the instance segmentation benchmark on the Cityscapes dataset compared to Mask R-CNN, which we are basing our work on. Yang et al. (Yang et al., 2016) propose an encoder - decoder architecture which predicts object contours. For training on the MS COCO dataset (Lin et al., 2014) the coarse polygon ground-truth edges are refined to follow the object contours more closely by applying a dense conditional random field (Krähenbühl and Koltun, 2011) or applying graph cut (Boykov and Jolly, 2001). This is particularly necessary in this case since the model predicts high resolution edges. Since the instance masks by Mask R-CNN only have a low resolution of e.g. 28 x 28, we argue that the preprocessing is unnecessary in this case, since most details are lost at this scale. It was also not used by the original Mask R-CNN Paper by He et al. (2017). Kirillov et al. (Kirillov et al., 2017) propose a model combination which predicts both a semantic segmentation and edges. The output of the edge score is used to compute super pixels which, alongside the predictions and the edge score, are used to solve a Multi Cut problem (Chopra and Rao, 1993) which predicts instances. This work is different from our experiments since we first find instance masks and then compare the predicted and ground-truth edges of these. In addition, this work operates on high resolution images, which complicates training, since the ground-truth edges only occupy a very small region of the overall area. The authors therefore rescale the cost function for underrepresented classes. This step is less relevant for the Edge Agreement Head applied to Mask R-CNN since the masks are predicted based on a proposed bounding box and because the class and mask prediction are decoupled in the Mask R-CNN architecture (He et al., 2017).

**Figure 1:** Overview of different example masks to illustrate the effect of the Edge Agreement Loss. *GT* corresponds to the ground truth and *I* to *IV* represent four example mask predictions which demonstrate early-stage predictions of the Mask R-CNN during training.

## 3. Edge Agreement Loss

When training a Mask R-CNN for instance segmentation one often observes incomplete or poor masks, especially during early training steps. Furthermore, the masks often do not follow the real object boundaries. Possible mistakes such as missing parts or oversegmentation are illustrated in Figure 1.

To reduce this problem, we draw our inspiration from how a human would perform instance segmentation: instead of immediately assigning parts of the image to specific objects one often identifies at first the boundaries of the object and fills the enclosed area. To help the network perform the segmentation in an analogous way, i.e. show the importance of edges and boundaries of objects, we have constructed an auxiliary loss called Edge Agreement Loss $L_{Edge}$. It is defined as the $L^p$ loss between the edges in the predicted mask and the ground-truth mask. The total loss $L_{Total}$ consists of the original Mask R-CNN loss $L_{MRNN}$ (eq. 1) and the new Edge Agreement Loss $L_{Edge}$ which are summed. To compute this new loss, the first step is to identify the edges in the predicted and the ground-truth mask.

### 3.1. Edge Detection

In detail, we examined edge detection filters which can be described as a convolution with a $3 \times 3$ kernel, such as the well-known *Sobel* and *Laplacian* filters.

The Sobel filters (Sobel and Feldman, 1973) are two-dimensional filters to detect edges. As the filters describe a first-order gradient operation they are rotation-dependent. There are two filters

$$S_x = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}, \quad S_y = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix} \quad (2)$$

which describe the horizontal and the vertical gradient respectively. An edge in the image corresponds to a high absolute response along the filter's direction. In the following the concatenation of both filters into a $3 \times 3 \times 2$ dimensional tensor is referenced as the Sobel filter $S$.

The Laplacian filter is a discretization of the two-dimensional Laplacian operator (i.e. the second derivative). The filter

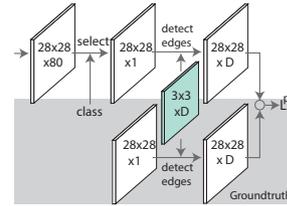$$L = \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad (3)$$

is the direct result of a finite-difference approximation of the derivative (Forsyth and Ponce, 2002). The operator is known

to be rotation invariant which means that it can detect edges in both x and y direction. As it is a second-order operator, an edge in the image corresponds to a zero-crossing, rather than a strong filter response. By including the main- and anti-diagonal elements the filter can be made responsive to 45° angles
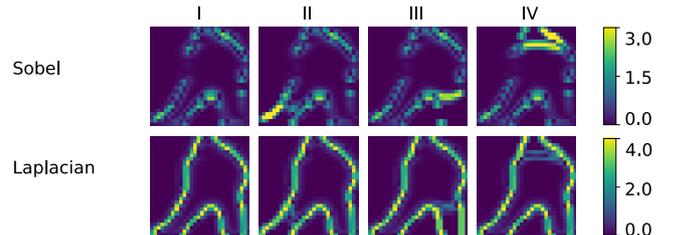
$$L = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{bmatrix}. \quad (4)$$

This is the Laplacian kernel ($L$) used in all further experiments.

In addition to these exemplary kernels we also used the Prewitt operator (Lipkin, 1970), Kayyali filter (Kawalec-Latała, 2014) and the Roberts operator (Roberts, 1963).



**Figure 2:** Edge Agreement Head: We extend the existing mask branch architecture. Of the $28 \times 28 \times 80$ dimensional output of the mask branch, the mask corresponding to the correct class is selected. The head computes a convolution of the selected mask and the ground-truth mask with the $3 \times 3 \times D$ dimensional edge detection filter (turquoise). Between these a $L^p$ loss is calculated, which results in the term $L_{Edge}$ (Best viewed in color).



**Figure 3:** $L^2$ errors for the four example predictions and the different methods. Each column *I* to *IV* corresponds to one of the examples in Figure 1. The first row shows the $L^2$ error based on the Sobel filter magnitude, while for the second row the Laplace filter is used.

### 3.2. Loss Construction

To calculate the final loss $L_{Edge}$ we propose an additional network head, called the Edge Agreement Head. It uses the predicted and the matched ground-truth masks as input, which are then convolved with a selection of edge detection filters. Afterwards, the difference between the predicted and ground-truth edge maps are determined. The entire procedure is illustrated in Figure 2 in the left half. For this task, we choose the set of $L^p$ loss functions. Mathematically they can be expressed as the $p$-th power of the generalized power mean $M_p$ of the absolute difference between the target $\hat{\mathbf{y}}$ and the prediction $\mathbf{y}$

$$L^p(\mathbf{y}, \hat{\mathbf{y}}) = M_p(|\mathbf{y} - \hat{\mathbf{y}}|)^p. \quad (5)$$

For $p = 2$ this equals the mean square error, commonly used in deep learning.

The edge agreement head can be calculated with only minimal additional computational and memory requirements. This means, that the method can be integrated in existing systems for training Mask R-CNN without requiring new or additional hardware.

## 4. Implementation Details

A mask size of $28 \times 28$ pixels and an image resolution of $1024 \times 800$ pixels are used. All training images are resized to this size preserving their aspect ratio. As the training images may have different aspect ratios, the remaining space of the image is zero padded. This method differs from the one used in the original Mask R-CNN implementation (He et al., 2017), where resizing is done such that the smallest side is 800 pixels and the largest is trimmed at 1000 pixels.

The ResNet (He et al., 2016) feature extractor is initialized with weights trained on ImageNet (Deng et al., 2009); all other weights (e.g. in the region proposal network) are initialized using Xavier initialization (Glorot and Bengio, 2010).

A similar training strategy to other Mask R-CNN work (He et al., 2017) is followed. We choose to train the network for 160k steps on the MS COCO 2017 train dataset with a batch size of 2 on a single GPU machine, while for Mask R-CNN an effective batch size of 16 was used. The training consists of three stages each lasting for 40k, 80k, 40k steps respectively: in the first stage only the Mask R-CNN branches and not the ResNet backbone are trained. Next, the prediction heads and parts of the backbone (starting at layer 4) are optimized. Finally, in the third stage, the entire model (backbone and heads) is trained together. For the first two training stages we use a learning rate of 0.001 and for the last one a decreased learning rate of 0.001/10. The optimization is done by SGD with momentum set to 0.9 and weight decay set to 0.0001.

## 5. Experiments

We perform our experiments using the implementation of Mask R-CNN by matterport (Abdulla, 2017), based on the KERAS framework (Chollet et al., 2015) with a TENSORFLOW backend (Abadi et al., 2015). Each training is carried out on a single GPU using either an NVIDIA Titan X or an NVIDIA GeForce GTX 1080 Ti.

We examine three aspects of the proposed Edge Agreement Head. At first, we inspect the influence of the edge detection filters on the training speed (section 5.1). In section 5.2, the different metrics of the $L^p$ family are used to examine the influence of the loss function's steepness on the training speed. Section 5.3 shows the impact weighting the Edge Agreement Loss has on the overall loss. We investigate the influence of the Edge Agreement Head with varied mask size in section 5.4. In section 5.5 we show the results on the metrics after longer training. Finally, in section 5.6 we elaborate on modifications to the

Edge Agreement Head which did not have a positive effect on the training.

For all experiments, we follow the same scheme: every network configuration examined is trained and evaluated three times. The resulting training curves and metrics displayed are the averaged values. Furthermore, to be able to compare all runs and to reduce the time required for all experiments we do not train the networks until they have converged, but only for a fixed and limited number of training steps. The only data augmentation used in all three steps are random horizontal flips.
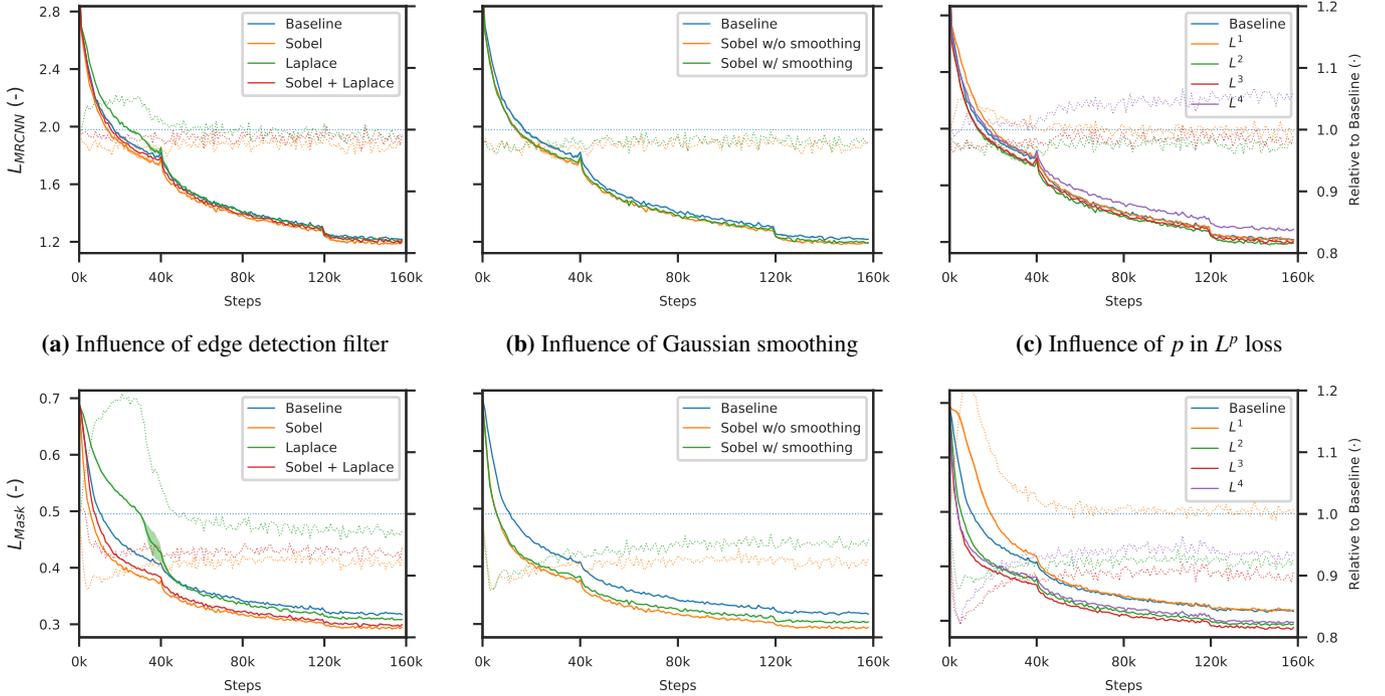
We present the COCO metrics for our experiments and compare them with the results obtained using a Mask R-CNN model without modifications (baseline) for every experiment conducted. A significant disadvantage of using the COCO Mask metrics is that they do not compare the ground-truth and predicted mask pixel per pixel since they only consider the area and the instance enclosing bounding box. As a result, the COCO Mask metrics are unsuitable to compare the improvement in the details of instance masks which are located on the inside of the mask and do not affect the extremities. We use the COCO metrics because of their dominance in other publications. The mask loss however, can be regarded as a better metric to compare the quality of the instance mask, because it is computed using a cross entropy between matched ground-truth and predicted mask.
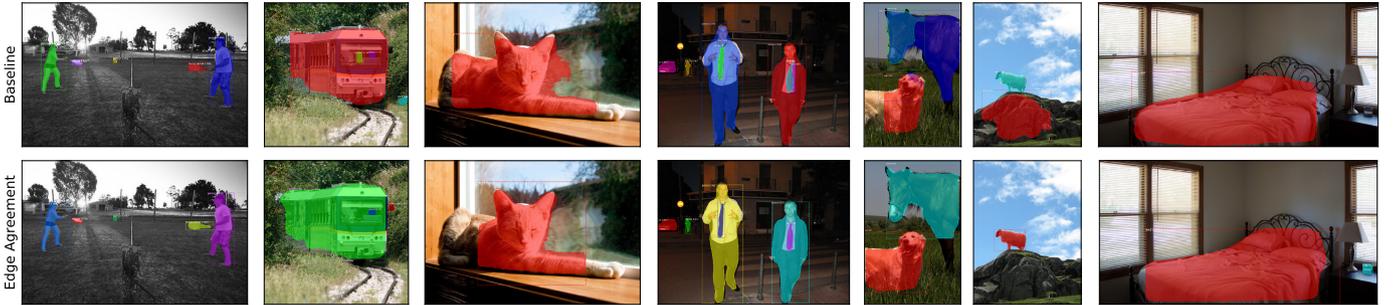
### 5.1. Influence of Filters

In the first experiment the choice of edge detection filters on the training speed and the mask quality is analyzed. The Edge Agreement Loss is computed using the $L^2$ loss. The mask loss $L_{Mask}$ and the original Mask R-CNN loss $L_{MRCNN}$ are displayed in Figure 4a. The graphs show that using the Sobel filter leads to a faster decrease of the $L_{MRCNN}$ and the $L_{Mask}$ loss. This is underlined well by plotting $L_{MRCNN}$ while using the Sobel filter relative to the baseline (dotted curves with respect to the right y-axis in Figure 4a) which demonstrates a consistent improvement at every training step.

To allow for a comparison between different edge detection filters, we analyzed their performance based on their impact on the other loss terms, as their Edge Agreement Loss magnitude varies depending on the filter. The results on the COCO Mask metrics are shown in Table 1. Note, that these results are lower than the results reported by (He et al., 2017), since we use a different batch size, input image resolution and implementation of Mask R-CNN ((Abdulla, 2017) not (Girshick et al., 2018)). Using the Sobel edge detection filter gave a 7% relative improvement on the AP score. Notably the Sobel filter resulted in a 12% relative improvement on the $AP_{75}$ and a 10% relative improvement on the $AP_S$ compared to the respective baseline scores. On average, the COCO metrics have been improved relatively by 8.1%. Using the Laplacian kernel showed only marginal improvements over the baseline. The combination of multiple filters, e.g. Sobel and Laplacian (S & L), showed no increase in performance.

A possible explanation for the superiority of the Sobel filter is its structure: as it consists of two filters, not only the strength of an edge along the x and y axis but also the edge's orientation

**Figure 4:** Comparison of different Edge Agreement Head configurations on the MS COCO dataset. The left y-axis corresponds to the absolute loss values (solid lines) and the right y-axis corresponds to the relative improvement compared to the baseline (dotted lines). The first row shows the original Mask R-CNN Loss $L_{MRCNN}$ while the second row shows the Mask Loss $L_{Mask}$. The first column illustrates the influence of different edge detectors used in the Edge Agreement Head, while the second demonstrates the influence of Gaussian smoothing when using a Sobel edge detection filter (see section 5.6). The last column compares the performance of different $L^p$ loss functions for the Edge Agreement Loss (Best viewed in color).



**Figure 5:** Comparison between masks predicted by Baseline Mask R-CNN and Mask R-CNN with Edge Agreement Head on the MS COCO dataset using a Sobel edge detection filter after 160k steps on images taken from the MS COCO 2017 dataset (Lin et al., 2014) (Best viewed in color).

can be used during the gradient descent to minimize the total loss. This additional information accelerates the training. Due to its similar structure, the Prewitt filter showed a comparable effect on the training.

A qualitative comparison between computed masks is shown in Figure 5. We observe that the models trained with Edge Agreement Loss tend to be less likely to propose bounding boxes which do not contain any object, therefore reducing false positives. This indicates that the features needed to minimize the Edge Agreement Loss are also useful to the Region Proposal Network.

### 5.2. Influence of the Choice p in $L^p$ loss

Next, the influence of the exponent $p$ in $L^p$ chosen for the Edge Agreement Loss is analyzed. For all previous experiments, an $L^2$ loss was used; now different values of $p \in \{1, 2, 3, 4\}$ are applied. As an increasing value of $p$ increases the steepness of the loss, falsely detected edges are penalized

more strongly. For this evaluation we used the Sobel edge detection filter without smoothing the ground truth.

The two losses $L_{MRCNN}$ and $L_{Mask}$ are displayed in Figure 4c. While a higher value for $p$ causes the mask loss to decrease, it also increases the overall loss. The metrics obtained in these experiments are listed in Table 2. Overall, choosing the $L^2$ loss appears to be the best choice, as it yields the best results on the COCO metrics.

### 5.3. Influence of Weighting Factor on Edge Agreement Loss

By choosing a higher value of $p$ for the $L^p$ Edge Agreement Loss, the loss becomes steeper and yields higher values for wrongly predicted masks. This increment also implies a higher relative importance of the Edge Agreement Loss compared to the other loss functions in the sum of the total loss which usually stay in the range [0, 1].

In these trainings we used the Sobel edge detection filter.

5

**Table 1:** Influence of the choice of edge detection filters on the instance segmentation mask COCO AP metrics on the MS COCO dataset after 160k steps. Higher is better.

| | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|
| Sobel | **20.2 ± 0.17** | **37.5 ± 0.37** | **20.0 ± 0.07** | **8.8 ± 0.27** | **21.9 ± 0.18** | **28.9 ± 0.30** |
| Prewitt | 20.0 ± 0.31 | **37.5 ± 0.25** | 19.6 ± 0.38 | 8.5 ± 0.45 | 21.6 ± 0.31 | 28.1 ± 0.61 |
| Kayyali | 19.7 ± 0.16 | 36.3 ± 0.25 | 19.5 ± 0.18 | 8.4 ± 0.32 | 21.3 ± 0.14 | 28.1 ± 0.46 |
| Roberts | 18.9 ± 0.31 | 36.4 ± 0.44 | 17.9 ± 0.28 | 7.9 ± 0.21 | 20.5 ± 0.41 | 26.7 ± 0.53 |
| Laplace | 19.4 ± 0.12 | 36.5 ± 0.19 | 18.9 ± 0.18 | 8.0 ± 0.22 | 21.0 ± 0.10 | 27.8 ± 0.11 |
| S & L | 20.0 ± 0.25 | 37.0 ± 0.36 | 19.6 ± 0.24 | 8.3 ± 0.03 | 21.7 ± 0.30 | 28.4 ± 0.47 |
| Baseline | 18.8 ± 0.14 | 36.5 ± 0.24 | 17.8 ± 0.13 | 8.0 ± 0.21 | 20.4 ± 0.29 | 26.6 ± 0.24 |

**Table 2:** Influence of the chosen $L^p$ loss on the instance segmentation mask AP COCO metrics on the MS COCO dataset after 160k steps. Higher is better.

| | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|
| $L^1$ | 19.5 ± 0.28 | 36.6 ± 0.41 | 18.9 ± 0.41 | 8.2 ± 0.30 | 21.0 ± 0.32 | 27.7 ± 0.5 |
| $L^2$ | **20.2 ± 0.17** | **37.5 ± 0.37** | 20.0 ± 0.07 | **8.8 ± 0.27** | **21.9 ± 0.18** | **28.9 ± 0.30** |
| $L^3$ | **20.2 ± 0.20** | 37.0 ± 0.41 | **20.1 ± 0.22** | 8.6 ± 0.14 | 21.8 ± 0.24 | 28.5 ± 0.53 |
| $L^4$ | 17.8 ± 0.13 | 33.5 ± 0.12 | 17.4 ± 0.13 | 7.6 ± 0.12 | 19.3 ± 0.23 | 24.7 ± 0.24 |
| Baseline | 18.8 ± 0.14 | 36.5 ± 0.24 | 17.8 ± 0.13 | 8.0 ± 0.21 | 20.4 ± 0.29 | 26.6 ± 0.24 |

To examine the influence of the relative importance of the new Edge Agreement Loss to the other losses, we include a factor $\alpha$ which scales the Edge Agreement Loss. We test its influence on the usage of the $L^2$ and $L^4$ losses to investigate the impact of the Mask Edge Loss on the total loss. For this comparison all trainings are performed only once and up to 120k steps instead of 160k steps, as already after 120k steps a clear trend has been recognizable.

Figure 6a shows the Mask R-CNN loss, while the Edge Agreement Loss is scaled by $\alpha \in \{0.5, 1, 8, 16\}$. The Mask R-CNN loss increases with higher weight factor, despite faster decreasing Mask Loss, the other loss terms remain higher. In fact, the $L^2$ loss with weight factor 1.0 already appears to be a good trade-off between enforcing better predicted masks and optimizing the other objectives of the network.

The $L^4$ loss yields high values compared to the $L^2$ loss. Therefore, we scale it by $\alpha = 1/16$, which is approximately the ratio of Edge Agreement Loss between using $L^2$ and $L^4$ in the first few steps. The training progression for the Mask R-CNN loss is shown in Figure 6b. Reducing the Edge Agreement Loss improves training significantly, making the loss stay below the Baseline for most of the steps.

### 5.4. Influence of Mask Size

We investigate the influence of the size of the predicted mask on the performance of the Edge Agreement Head (see Table 3). For this we compare the performance of models trained for the original mask size of 28×28 and for an increased size of 56×56. The models were trained following the same training schedule as outlined in section 4. For both mask sizes we observe a clear increase in performance when using the Edge Agreement Head as an auxiliary loss. However, we find that the overall performance of the model trained at mask size 56 × 56 is worse than the model trained with mask size 28 × 28. We hypothesize that this is the case, because predicting twice the resolution would require longer training.
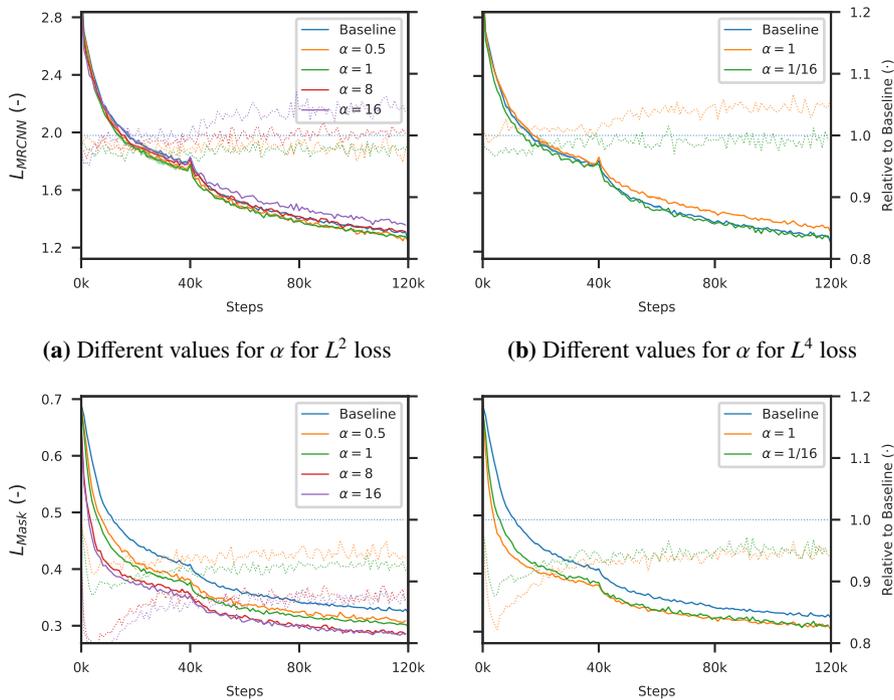
### 5.5. Longer training

To measure the effect of the Edge Agreement Head after longer training, we increased the number of steps previously used (320k and 640k steps rather than 160k). In this case the last step of the training schedule, in which all layers are trained, was extended from 40k steps to 200k and 520k steps (a total of 320k and 640k training steps respectively). The results are shown in Table 4.

All metrics improved as expected when trained for additional steps. Interestingly, in most metrics the Mask R-CNN model trained with Edge Agreement Head trained for 160k steps was not only superior to the baseline trained for 160k steps but also to the one trained 320k steps. No significant influence of the Edge Agreement Loss on losses other than the Mask Loss is observed. We notice that the difference in the Mask Loss between a baseline Mask R-CNN and one trained with Edge Agreement Loss remains constant with later training steps. This was contrary to our own intuition that the Edge Agreement Loss would primarily be helpful early in training. It was expected that the Mask Loss of the two models would approach each other, but this was not found to be the case. We conclude that the Edge Agreement Head is not only useful early on in training, but can guide the training even in later training steps and change the point of convergence.

It should be noted that our results on the MS COCO dataset (Lin et al., 2014) are significantly lower than the results reported by He et al. (2017). Firstly, we are not using the official implementation of Mask R-CNN made available in the Detectron (Girshick et al., 2018), but an independent implementation which reported lower results of their pretrained models (Abdulla, 2017). Secondly, we use a batch size of 2, while Mask R-CNN used an effective batch size of 16. We argue that this does not hurt the generality of our method.

### 5.6. Other experiments

The configuration of the Edge Agreement Loss we describe above was found to have the optimal impact on the training. We

**Figure 6:** Influence of the weighting factor $\alpha$ on the behavior of the Edge Agreement Loss on the MS COCO dataset. The first row displays again the original total loss of Mask R-CNN $L_{MRCNN}$ while the second row displays only the mask loss $L_{Mask}$. The first column shows the loss trajectory for different alpha values using the $L^2$ loss whereas the second column shows the influence on the $L^4$ loss. Best viewed in color

**(a)** Different values for $\alpha$ for $L^2$ loss

**(b)** Different values for $\alpha$ for $L^4$ loss

**Table 3:** Influence of the size of the predicted and ground-truth masks on the Edge Agreement Head. Shown are the instance segmentation mask AP COCO metrics on the MS COCO dataset.

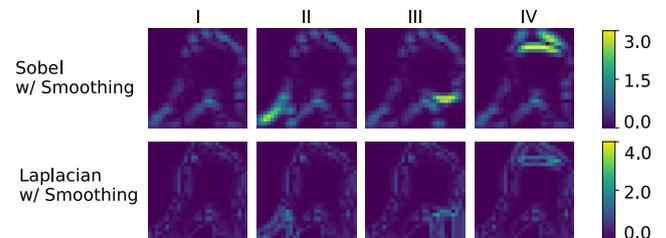| Mask shape | | **AP** | **AP$_{50}$** | **AP$_{75}$** | **AP$_S$** | **AP$_M$** | **AP$_L$** |
|---|---|---|---|---|---|---|---|
| $28 \times 28$ | Baseline | $18.8 \pm 0.14$ | $36.5 \pm 0.24$ | $17.8 \pm 0.13$ | $8.0 \pm 0.21$ | $20.4 \pm 0.29$ | $26.6 \pm 0.24$ |
| | Ours | $\mathbf{20.2 \pm 0.17}$ | $\mathbf{37.5 \pm 0.37}$ | $\mathbf{20.0 \pm 0.07}$ | $\mathbf{8.8 \pm 0.27}$ | $\mathbf{21.9 \pm 0.18}$ | $\mathbf{28.9 \pm 0.30}$ |
| $56 \times 56$ | Baseline | $18.0 \pm 0.23$ | $35.0 \pm 0.30$ | $17.0 \pm 0.37$ | $7.6 \pm 0.24$ | $19.3 \pm 0.28$ | $25.5 \pm 0.31$ |
| | Ours | $19.3 \pm 0.03$ | $36.0 \pm 0.07$ | $19.0 \pm 0.07$ | $8.1 \pm 0.13$ | $21.0 \pm 0.07$ | $27.6 \pm 0.19$ |

tried a variety of modifications which showed either no effect or had a negative impact on the training.

### 5.6.1. Smoothing of ground-truth or predicted Masks

Figure 3 illustrates how the $L^2$ loss between the edge maps, calculated as mentioned above, does not only contain important information but even possibly distracting information: the images in the rows for Sobel and Laplace depict a high error rate which is not limited to areas in which the person has not been segmented but also occurs around the entire boundary. This is mainly caused by the fact, that the ground-truth mask is binary and the mask branch's output is continuous and shows often smooth transitions.

To overcome this problem, we add an additional step in our branch which performs Gaussian smoothing on the ground-truth mask, yielding a smooth version of the binary ground truth. For this we use an approximate $3 \times 3$ Gaussian kernel. The calculated $L^p$ distance using this head proposal is shown in Figure 7. Notably, the loss calculated on the smoothed ground truth focuses particularly on areas with missing parts, while the Edge Agreement Loss on the default ground truth has a higher value on the overall mask boundary.

Contrary to expectation, the Sobel filter was more effective when used without smoothing the ground truth, as the Mask R-CNN training loss fell faster and was lower in this case, as shown in Figure 4b. Particularly the Mask Loss during training



**Figure 7:** $L^2$ errors for the four example predictions and the different filters calculated on the Gaussian smoothed ground truth. Each column $I$ to $IV$ corresponds to one of the examples in Figure 1. The first row shows the $L^2$ error based on the Sobel filter magnitude, while for the second row the Laplace filter is used. While the losses calculated on the default ground truth (see Figure 3) do not respond strongly to missing areas, the losses on the smoothed ground truth are particularly pronounced in these areas.

is lower.

The results obtained contradict the theoretical considerations of the possible benefit to smoothing the ground truth. The smoothing of the ground truth was designed to ignore minor mistakes at the boundary of an almost perfectly predicted mask, but only focusing on major mistakes. Apparently, the network does not only profit from highlighting the most crucial mistakes in the predicted masks, but rather from all mistakes done.

We investigated whether smoothing both the ground truth and the predicted mask or only the predicted mask would help the network during training. The reasoning for this was that it

**Table 4:** Comparison of the instance segmentation mask AP COCO metrics on the MS COCO dataset of our best performing model with the baseline after an extended training duration. The best performing model uses the Edge Agreement Head with Sobel edge detection filter and $L^2$ Edge Agreement Loss.

| | **AP** | **AP$_{50}$** | **AP$_{75}$** | **AP$_S$** | **AP$_M$** | **AP$_L$** |
|---|---|---|---|---|---|---|
| Ours 160k steps | $20.2 \pm 0.17$ | $37.5 \pm 0.37$ | $20.0 \pm 0.07$ | $8.8 \pm 0.27$ | $21.9 \pm 0.18$ | $28.9 \pm 0.30$ |
| Ours 320k steps | 21.3 | 38.7 | 21.1 | 8.9 | 23.2 | 30.0 |
| Ours 640k steps | **22.7** | **41.0** | **23.1** | **10.2** | **24.6** | **32.0** |
| Baseline 160k steps | $18.8 \pm 0.14$ | $36.5 \pm 0.24$ | $17.8 \pm 0.13$ | $8.0 \pm 0.21$ | $20.4 \pm 0.29$ | $26.6 \pm 0.24$ |
| Baseline 320k steps | 20.0 | 38.5 | 19.1 | 8.6 | 21.6 | 28.2 |
| Baseline 640k steps | 21.5 | 40.5 | 20.8 | 9.0 | 22.9 | 30.7 |

could be beneficial to penalize the network less for pixel accurate mask and more for the general shape. Since the instance boundaries become much wider due to the smoothing, the Edge Agreement Loss becomes less sensitive to small spatial displacements. However, we found this modification to have a negative impact on the training.

### 5.6.2. Balancing Losses

As discussed in section 5.3, the magnitude of the Edge Agreement Loss appears to have a high influence on the $L_{MRCNN}$ loss. In an attempt to balance the loss terms we tried homoscedastic task uncertainty as proposed by Kendall et al. (2018). Our approach was to weigh all the loss terms including the Edge Agreement Loss. However, the results were consistently worse than the baseline and therefore not included in this paper.

### 5.6.3. Alternative Edge Loss Definitions

Furthermore, we tried to weigh the cross entropy mask loss $L_{Mask}$ with the Edge Agreement Loss. Two different formulations for this weighted cross entropy loss were tried out, which can be expressed as

$$L_{Edge} = L_{Mask-PW} \cdot L_{Edge-PW}$$
$$\text{or}$$
$$L_{Edge} = L_{Mask-PW} \cdot \exp\left(L_{Edge-PW}/4\right),$$

using $L_{Mask-pw}$ and $L_{Edge-pw}$ to denote the pixel-wise Mask Loss and pixel-wise Edge Agreement Loss respectively. For both formulations the results were identical with the more concise formulation of the Edge Agreement Loss that we used in the rest of the paper.

In addition, when using the Sobel filter, we did not solely consider the horizontal and vertical image gradient but also the gradient's magnitude for calculating the $L^p$ Edge Agreement Loss. No improvement compared to not including the magnitude was found. Therefore, it was not used in the rest of the paper.

### 5.7. Cityscapes

To verify that our findings could be reproduced on a different dataset, we trained our models on the Cityscapes dataset (Cordts et al., 2016) with the Edge Agreement Head with Mask shape $28 \times 28$ and $56 \times 56$ pixels (see Table 5). In contrast to MS COCO, the annotations in Cityscapes have much finer details. We followed the training schedule of the authors of

Mask R-CNN (He et al., 2017), but instead of using an effective batch size of 8 we used a reduced size of 4. The findings are easily compared to Table 3, since we see very similar trends. Training with the Edge Agreement Head is demonstrated to be beneficial since it consistently outperforms the respective baseline. Increasing the predicted mask size leads to an overall reduction in the accuracy of the predicted mask for the baseline experiments, but this could potentially be remedied by longer training, to account for the higher number of parameters in the additional layer of the mask branch.

## 6. Conclusion

In this paper we have analyzed the behavior of Mask R-CNN networks during early training steps. By inspecting the predicted masks of the mask branch, we recognized that these often have blurry boundaries which do not follow sharp and fine contours of the original masks. To reduce this symptom, we successfully introduced a parameter free network head, the Edge Agreement Head. This head uses classical edge detection filters applied on the instance masks to calculate a $L^p$ loss between the predicted and ground-truth mask contours.

By including the new Edge Agreement Loss in the training, we achieved a relative performance increment of 8.1% averaged over all the MS COCO metrics after a fixed number of 160k training steps.

The ablation studies performed showed that the Sobel filter yields a better performance than the Laplace filter. Beyond expectations, the proposed smoothing of the ground-truth mask did not improve but hinder the performance. Out of all losses examined the often-used $L^2$ loss performs the best.

When trained longer, the difference in Mask Loss between a baseline Mask R-CNN and one with Edge Agreement Head persists, demonstrating the effectiveness of the additional loss not only early during training but also during later steps.

Finally, we demonstrated that the Edge Agreement Head is beneficial on Cityscapes, a dataset with much finer ground-truth masks.

## 7. Future work

The idea to enforce edge agreement in predicted semantic segmentation could be applied to scene segmentation for example on the DeepLab architecture (Chen et al., 2018a) or the U-Net architecture (Ronneberger et al., 2015; Kohl et al., 2018). Monocular depth estimation could also potentially be enhanced

**Table 5:** Influence of the size of the predicted and ground-truth masks on the Edge Agreement Head with Sobel and an $L^2$ loss. Shown are the instance segmentation mask AP COCO metrics on the Cityscapes dataset.

| Mask shape | | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|
| $28 \times 28$ | Baseline | $15.6 \pm 0.84$ | $33.5 \pm 1.8$ | $12.43 \pm 1.03$ | $2.2 \pm 0.25$ | $11.2 \pm 0.43$ | $25.6 \pm 1.21$ |
| | Ours | $17.7 \pm 0.49$ | $\mathbf{36.1 \pm 0.78}$ | $14.5 \pm 0.74$ | $\mathbf{3.0 \pm 0.10}$ | $12.5 \pm 0.81$ | $\mathbf{33.3 \pm 4.54}$ |
| $56 \times 56$ | Baseline | $14.9 \pm 0.40$ | $31.5 \pm 0.22$ | $11.9 \pm 0.63$ | $2.1 \pm 0.11$ | $10.2 \pm 0.97$ | $24.4 \pm 0.82$ |
| | Ours | $\mathbf{18.0 \pm 1.09}$ | $35.0 \pm 2.65$ | $\mathbf{15.05 \pm 0.99}$ | $1.8 \pm 0.04$ | $\mathbf{12.7 \pm 0.64}$ | $31.5 \pm 0.297$ |

by encouraging the predicted depth map to have comparable gradients to the ground-truth depth map image gradients.

Furthermore, balancing the different individual losses contained in the total loss by introducing new scaling variables might be a necessary step to further increase the training speed. Instead of introducing new static hyperparameters for the multi-task loss one could modify the gradients like Chen et al. (2018b).

As the Edge Agreement Loss accelerates the training of the Mask Head, it enables Mask R-CNN to be used more easily with sparse labels for object instance masks. This allows new training strategies for new datasets, e.g. one could mix a few hand segmented frames with datasets containing only object bounding boxes, such as PASCAL VOC (Everingham et al., 2010) or the more recent Open Images dataset (Krasin et al., 2017).

## Acknowledgments

## References

## References

Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.

Abdulla, W. (2017). Mask R-CNN for object detection and instance segmentation on Keras and TensorFlow.

Bai, M. and Urtasun, R. (2017). Deep watershed transform for instance segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Bertasius, G., Shi, J., and Torresani, L. (2015). Deepedge: A multi-scale bifurcated deep network for top-down contour detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Borgefors, G. (1986). Distance transformations in digital images. *Computer vision, graphics, and image processing*, 34(3):344–371.

Boykov, Y. Y. and Jolly, M.-P. (2001). Interactive graph cuts for optimal boundary & region segmentation of objects in nd images. In *Computer Vision, 2001. ICCV 2001. Proceedings. Eighth IEEE International Conference on*, volume 1, pages 105–112. IEEE.

Caruana, R. (1997). Multitask learning. *Mach. Learn.*, 28(1):41–75.

Chen, L., Papandreou, G., Kokkinos, I., Murphy, K., and Yuille, A. L. (2018a). Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848.

Chen, L.-C., Barron, J. T., Papandreou, G., Murphy, K., and Yuille, A. L. (2016). Semantic image segmentation with task-specific edge detection using cnns and a discriminatively trained domain transform. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Chen, Z., Badrinarayanan, V., Lee, C.-Y., and Rabinovich, A. (2018b). Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *ICML*.

Chollet, F. et al. (2015). Keras. https://keras.io.

Chopra, S. and Rao, M. R. (1993). The partition problem. *Mathematical Programming*, 59(1-3):87–115.

Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM.

Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., and Schiele, B. (2016). The cityscapes dataset for semantic urban scene understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee.

Everingham, M., Van Gool, L., Williams, C. K. I., Winn, J., and Zisserman, A. (2010). The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338.

Forsyth, D. A. and Ponce, J. (2002). *Computer Vision: A Modern Approach*. Prentice Hall Professional Technical Reference.

Girshick, R. (2015). Fast r-cnn. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1440–1448.

Girshick, R., Radosavovic, I., Gkioxari, G., Dollár, P., and He, K. (2018). Detectron. https://github.com/facebookresearch/detectron.

Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256.

Godard, C., Mac Aodha, O., and Brostow, G. J. (2017). Unsupervised monocular depth estimation with left-right consistency. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Hayder, Z., He, X., and Salzmann, M. (2017). Boundary-aware instance segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

He, K., Gkioxari, G., Dollar, P., and Girshick, R. (2017). Mask r-cnn. In *The IEEE International Conference on Computer Vision (ICCV)*.

He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Kawalec-Latała, E. (2014). Edge detection on images of pseudoimpedance section supported by context and adaptive transformation model images. *Studia Geotechnica et Mechanica*, 36(1):29–36.

Kendall, A., Gal, Y., and Cipolla, R. (2018). Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Kirillov, A., Levinkov, E., Andres, B., Savchynskyy, B., and Rother, C. (2017). Instancecut: from edges to instances with multicut. In *CVPR*, volume 3, page 9.

Kohl, S., Romera-Paredes, B., Meyer, C., De Fauw, J., Ledsam, J. R., Maier-Hein, K., Eslami, S. A., Rezende, D. J., and Ronneberger, O. (2018). A probabilistic u-net for segmentation of ambiguous images. In *Advances in*

*Neural Information Processing Systems*, pages 6965–6975.

Konishi, S., Yuille, A. L., Coughlan, J. M., and Zhu, S. C. (2003). Statistical edge detection: Learning and evaluating edge cues. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(1):57–74.

Krähenbühl, P. and Koltun, V. (2011). Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in neural information processing systems*, pages 109–117.

Krasin, I., Duerig, T., Alldrin, N., Ferrari, V., Abu-El-Haija, S., Kuznetsova, A., Rom, H., Uijlings, J., Popov, S., Kamali, S., Malloci, M., Pont-Tuset, J., Veit, A., Belongie, S., Gomes, V., Gupta, A., Sun, C., Chechik, G., Cai, D., Feng, Z., Narayanan, D., and Murphy, K. (2017). Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from https://storage.googleapis.com/openimages/web/index.html*.

Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc.

Li, Y., Qi, H., Dai, J., Ji, X., and Wei, Y. (2017). Fully convolutional instance-aware semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Liebel, L. and Körner, M. (2018). Auxiliary tasks in multi-task learning. *arXiv preprint arXiv:1805.06334*.

Lin, T.-Y., Dollar, P., Girshick, R., He, K., Hariharan, B., and Belongie, S. (2017). Feature pyramid networks for object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. (2014). Microsoft coco: Common objects in context. In Fleet, D., Pajdla, T., Schiele, B., and Tuytelaars, T., editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham. Springer International Publishing.

Lipkin, B. S. (1970). *Picture processing and psychopictorics*. Elsevier.

Liu, S., Qi, L., Qin, H., Shi, J., and Jia, J. (2018). Path aggregation network for instance segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., and Berg, A. C. (2016). Ssd: Single shot multibox detector. In *European Conference on Computer Vision (ECCV)*, pages 21–37. Springer.

Long, J., Shelhamer, E., and Darrell, T. (2015). Fully convolutional networks for semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Marmanis, D., Schindler, K., Wegner, J. D., Galliani, S., Datcu, M., and Stilla, U. (2018). Classification with an edge: Improving semantic image segmentation with boundary detection. *ISPRS Journal of Photogrammetry and Remote Sensing*, 135:158–172.

Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Ren, S., He, K., Girshick, R., and Sun, J. (2017). Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149.

Roberts, L. G. (1963). *Machine perception of three-dimensional solids*. PhD thesis, Massachusetts Institute of Technology.

Ronneberger, O., Fischer, P., and Brox, T. (2015). U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer.

Ruder, S. (2017). An overview of multi-task learning in deep neural networks. *CoRR*, abs/1706.05098.

Shen, W., Wang, X., Wang, Y., Bai, X., and Zhang, Z. (2015). Deepcontour: A deep convolutional feature learned by positive-sharing loss for contour detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Simonyan, K. and Zisserman, A. (2015). Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*.

Sobel, I. and Feldman, G. (1973). A 3x3 isotropic gradient operator for image processing. *in Hart, P. E. & Duda R. O. Pattern Classification and Scene Analysis*.

Szegedy, C., Ioffe, S., Vanhoucke, V., and Alemi, A. (2017). Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI Conference on Artificial Intelligence*.

Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. (2015). Going deeper with convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Xie, S. and Tu, Z. (2015). Holistically-nested edge detection. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1395–1403.

Yang, J., Price, B., Cohen, S., Lee, H., and Yang, M.-H. (2016). Object contour detection with a fully convolutional encoder-decoder network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 193–202.

Zhang, Z., Luo, P., Loy, C. C., and Tang, X. (2014). Facial landmark detection by deep multi-task learning. In *European Conference on Computer Vision (ECCV)*, pages 94–108. Springer.

## Authors

**Roland S. Zimmermann** received his B.Sc. degree with distinction in Physics from Georg-August University of Göttingen, Göttingen, Germany in 2017 where he is currently pursuing his M.Sc. in collaboration with the University of Tübingen, Tübingen, Germany. His research interests lie in the areas of computer vision and (adversarial) robustness of neural networks.

**Julien N. Siems** received his B.Sc. degree in Computer Science from TU Dresden, Dresden, Germany in 2017. He is currently pursuing his M.Sc. at the Albert Ludwig University of Freiburg, Freiburg, Germany. His research interests lie in the areas of computer vision and neural architecture search.