

Semi-Supervised Semantic Segmentation with High- and Low-level Consistency

Sudhanshu Mittal, Maxim Tatarchenko, and Thomas Brox

Abstract—The ability to understand visual information from limited labeled data is an important aspect of machine learning. While image-level classification has been extensively studied in a semi-supervised setting, dense pixel-level classification with limited data has only drawn attention recently. In this work, we propose an approach for semi-supervised semantic segmentation that learns from limited pixel-wise annotated samples while exploiting additional annotation-free images. It uses two network branches that link semi-supervised classification with semi-supervised segmentation including self-training. The dual-branch approach reduces both the low-level and the high-level artifacts typical when training with few labels. The approach attains significant improvement over existing methods, especially when trained with very few labeled samples. On several standard benchmarks - PASCAL VOC 2012, PASCAL-Context, and Cityscapes - the approach achieves new state-of-the-art in semi-supervised learning.

Index Terms—Computer Vision, Semi-supervised Learning, Semantic Segmentation, Generative Adversarial Networks.



1 INTRODUCTION

SEMANTIC segmentation is one of the key computer vision tasks important in various applications including autonomous driving, medical-imaging and robotics. Lately, Deep Convolutional Neural Networks [19] have demonstrated great results on this task for different datasets [4], [5], [31], [40]. However, this success usually comes at the cost of collecting dense pixel-wise annotations - a cumbersome process that involves much manual effort.

Attempting to alleviate the problem, several methods exploit weaker forms of supervision: image-level labels [1], [32], [38], bounding boxes [7], [29], or scribbles [20], [36]. Only two previous works [15], [34] have considered true semi-supervised learning for semantic segmentation, which requires having a small subset of fully-labeled samples along with a larger set of completely annotation-free images.

In this work, we propose a dual-branch method for semi-supervised semantic segmentation which can effectively learn from annotation-free samples given a very small set of fully-annotated samples. Our design is based on the observation that CNNs trained on limited data are subject to two typical modes of failure; see Figure 1(c-d). The first one appears as inaccuracy in low-level details, such as wrong object shapes, inaccurate boundaries, and incoherent surfaces. The second one is the misinterpretation of high-level information, which leads to assigning large image regions to wrong classes.

The two network branches are designed to separately address those two types of artifacts. To deal with low-level errors, we propose an improved GAN-based model, where the segmentation network acts as a generator. It is trained together with a discriminator that classifies between generated and ground-truth segmentation maps. Instead of

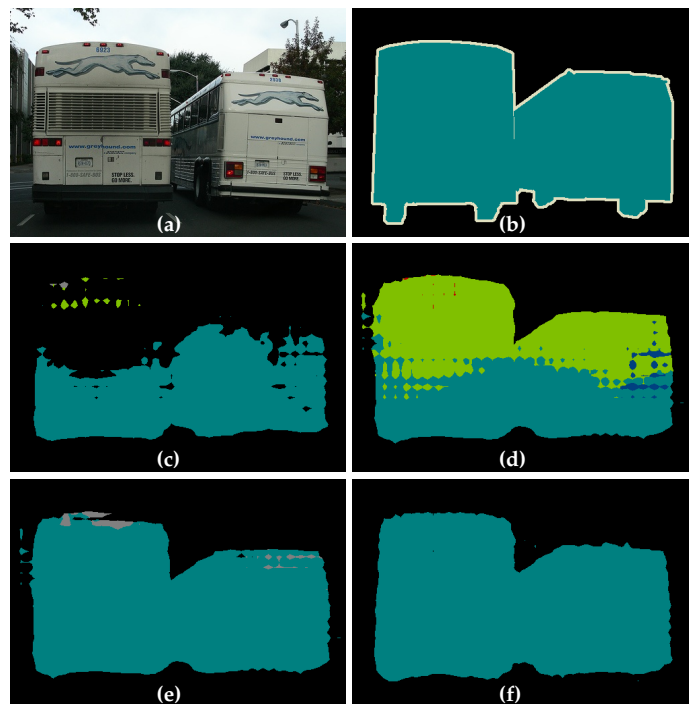


Fig. 1. An image from the PASCAL VOC dataset (a) and its ground-truth segmentation mask (b). Prediction (c) is obtained with supervised training on 5% labeled samples. Using the other 95% unlabeled images, our GAN-based branch improves the shape estimation (d). The second branch adds high-level consistency by removing false positives (e). (f) shows the output when training on 100% pixel-wise labeled samples.

using the original GAN loss, we propose to use the feature matching loss introduced by Salimans *et al.* [33]. Moreover, we introduce the self-training procedure based on the discriminator score which improves the final performance via leveraging high-quality generator predictions as fully labeled samples.

For the second type of artifacts, we propose a semi-

• The authors are with the Computer Science Department at the University of Freiburg, Freiburg im Breisgau, Germany.

E-mail: {mittal, tatarchm, brox}@cs.uni-freiburg.de

supervised multi-label classification branch which decides on the classes present in the image and thus aids the segmentation network to make globally consistent decisions. To utilize extra image-level information from unlabeled images, we leverage the success of ensemble-based semi-supervised classification (SSL) methods [18], [37].

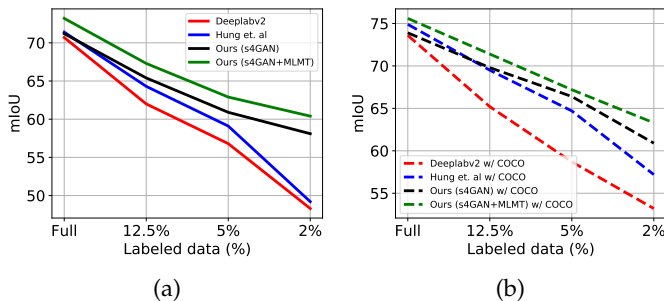


Fig. 2. **Semi-supervised Semantic Segmentation:** The proposed semi-supervised learning (SSL) approach improves over the baselines even when only little labeled data is available using unlabeled data, shows considerable improvement especially with less than 5% labeled samples. Performance is shown on the PASCAL VOC dataset without (a) and with (b) COCO pre-training.

The two branches act in a complementary manner and successfully fix both low-level and high-level errors; see Fig. 1 for a typical example. We demonstrate the effectiveness of our approach on different amounts of labeled data across a range of popular semantic segmentation datasets: PASCAL VOC 2012 [9], PASCAL-Context [27] and Cityscapes [6]. We consistently achieve the best results compared to existing methods and define the new state of the art in semi-supervised semantic segmentation. Our approach proves particularly efficient when only very few training samples are available: with as little as 2% labeled data we report an 11% performance improvement over the state of the art (see Figure 2).

We further show that the approach can easily make use of extra image-level weak annotations when those are available. It compares favorably to the existing methods operating in the same setting. The source code of this paper is available ¹.

2 RELATED WORK

Weakly-supervised and Semi-supervised Segmentation.

To reduce annotation effort, most existing approaches rely on weakly- and semi-supervised training schemes which use weak labels from the whole dataset like image-level class labels [29], [39], bounding boxes [7], [16], [29] or scribbles [20], [36], where semi-supervised schemes [16], [29], [39] additionally use a few pixel-wise segmentation labels.

Only two recent works [15], [34] consider true semi-supervised learning, i.e., they improve semantic segmentation with completely annotation-free images. These methods, like ours, utilize a GAN-based model. However, both approaches use the GAN in a different manner. Souly *et al.* [34] use the GAN to generate additional images to enhance the features learned by the segmentation network. They

further extend their semi-supervised method by generating additional class-conditional images.

Most related to ours is the work by Hung *et al.* [15]. They also propose a GAN-based design which enables learning from unlabeled samples. However, our framework is substantially different in many details. Hung *et al.* use an FCN-based [23] discriminator which yields a dense probabilistic map for each pixel, whereas we propose an image-wise discriminator. In contrast to the two-stage training process of [23], we propose an automatic integration of self-training based on the GAN training dynamics. Moreover, we propose to use a feature matching loss, which is crucial for the stability of GAN training, especially when only few labeled samples are available. Finally, we add a semi-supervised multi-label classification branch for resolving high-level inconsistencies.

Also Luc *et al.* [24] share some common ground with our work, although their work does not comprise semi-supervised learning. In their case the GAN replaces CRF-post-processing which enhance low-level consistency in the segmentation maps. Luc *et al.* [24] optimize the original GAN loss to encourage predicted segmentation maps to be similar to the ground-truth maps and show that it improves the performance in a fully-supervised setting.

Semi-supervised Classification. In contrast to segmentation, many semi-supervised methods exist for image classification [2], [18], [26], [33], [37]. Oliver *et al.* [28], however, criticize that most of the work lacks realistic evaluation to address real-world conditions. They propose a new experimental methodology closer to the real-world settings. We find that consistency-based semi-supervised classification methods [2], [37] show improvement over the supervised baseline while satisfying at least two procedures mentioned by Oliver *et al.* [28]. Firstly, those methods show improvement over the supervised setting while using a high-quality supervised baseline. Secondly, they can improve upon the pre-trained network using unlabeled data. We use the Mean-Teacher method [37] in our approach.

Network Fusion. The approach to fuse spatial and class information by channel-wise selection is inspired by some recent works in other domains. Hu *et al.* [14] proposed SE-Net for image classification, which learns to combine spatial and channel-wise information by calibrating channel-wise feature maps. Following SE-Net, Zhang *et al.* [40] proposed to incorporate class information in semantic segmentation to highlight class-dependent feature maps. Multiple works [13], [38], [39] have explored the usage of classification methods, both in a shared and a decoupled manner to constructively use class information for semi- and weakly supervised semantic segmentation. In this work, we use a decoupled approach with late fusion of spatial and class information to remove false positive class channels.

3 METHOD

We propose a two-branch approach to the task of semi-supervised semantic segmentation as shown in Figure 3. The lower branch predicts pixel-wise class labels and is referred to as the *Semi-Supervised Semantic Segmentation GAN* (s4GAN). The upper branch performs image-level classification and is denoted as the *Multi-Label Mean Teacher* (MLMT).

1. Source code: <https://github.com/sud0301/semisup-semseg>

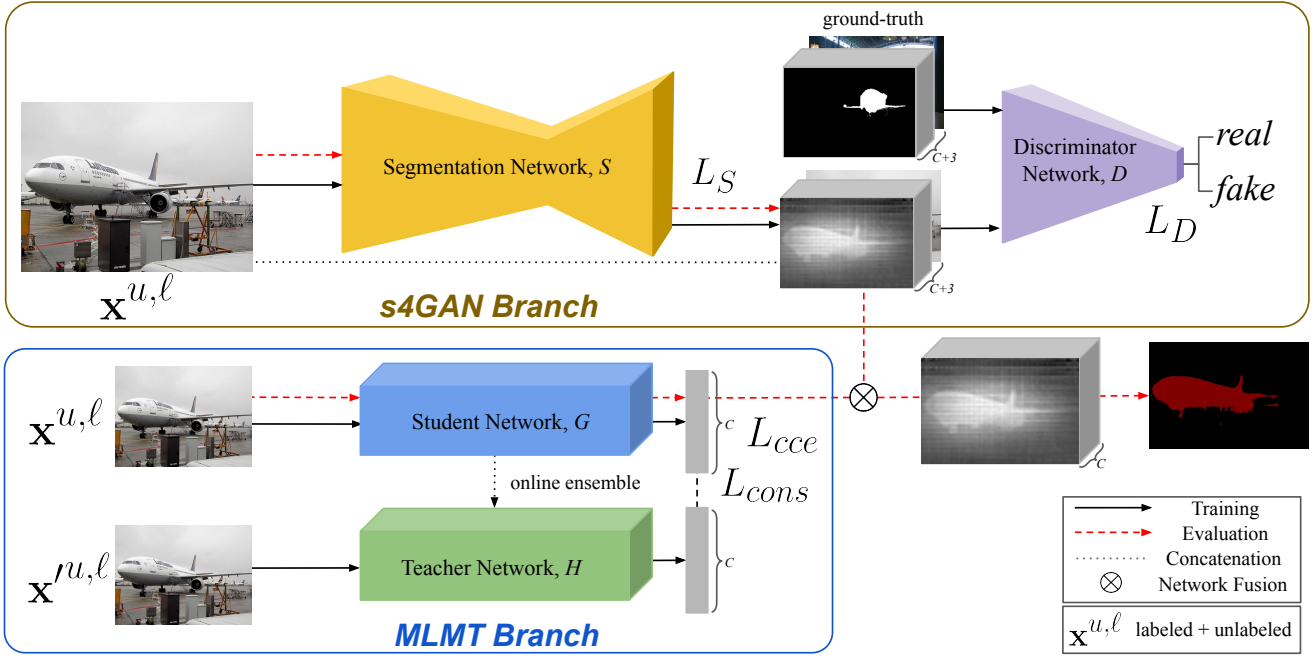


Fig. 3. Overview of our proposed semi-supervised segmentation approach. The s4GAN branch is a GAN-based model which improves the low-level details in the segmentation prediction. The MLMT branch performs semi-supervised multi-label classification to exploit class-level information for removing false-positive predictions from the segmentation map.

The core of the s4GAN branch is a standard segmentation network for generating per-pixel class labels given the input image. We combine conventional supervised training with adversarial training, which allows leveraging unlabeled data to improve the prediction quality. The segmentation network acts as a generator and is trained together with a discriminator responsible for distinguishing the ground truth segmentation maps from the generated ones. We additionally treat the outputs of the discriminator as a quality measure and use it to identify the best predictions which are further exploited for self-training.

The MLMT branch predicts image-level class labels used to filter the s4GAN outputs. Its core is a Mean Teacher classifier, which effectively removes false positive predictions of the segmentation network. The contributions of the two branches are complementary to each other. Their outputs are combined to produce the final result.

Notations: Our dataset \mathcal{D} is split into the labeled part $\mathcal{D}^\ell = \{\mathbf{x}^\ell, \mathbf{y}^\ell\}$ and the unlabeled part $\mathcal{D}^u = \{\mathbf{x}^u\}$, where \mathbf{x} are the input images and \mathbf{y} are the pixel-wise segmentation labels.

3.1 s4GAN for Semantic Segmentation

In our s4GAN model, the segmentation network S acts as a generator network that takes image \mathbf{x} as input and predicts C segmentation maps, one for each class. The discriminator D gets the concatenated input of the original image and its corresponding predicted segmentation. Its task is to match the distribution statistics of the predicted and the real segmentation maps.

3.1.1 Training S

The segmentation network S is trained with loss L_S , which is a combination of three losses: the standard cross-entropy

loss, the feature matching loss, and the self-training loss.

Cross-entropy loss. This is a standard supervised pixel-wise cross-entropy loss term L_{ce} . The loss for the output $S(\mathbf{x})$ of size $H \times W \times C$ is evaluated only for the labeled samples \mathbf{x}^ℓ :

$$L_{ce} = - \sum_{h,w,c} \mathbf{y}^\ell(h, w, c) \log S(\mathbf{x}^\ell)(h, w, c), \quad (1)$$

where \mathbf{y}^ℓ is the ground-truth segmentation mask.

Feature matching loss. The feature matching loss L_{fm} [33] aims to minimize the mean discrepancy between the feature statistics of the predicted, $S(\mathbf{x}^u)$ and the ground-truth segmentation maps, \mathbf{y}^ℓ :

$$L_{fm} = \left\| \mathbb{E}_{(\mathbf{x}^\ell, \mathbf{y}^\ell) \sim \mathcal{D}^\ell} [D_k(\mathbf{y}^\ell \oplus \mathbf{x}^\ell)] - \mathbb{E}_{\mathbf{x}^u \sim \mathcal{D}^u} [D_k(S(\mathbf{x}^u) \oplus \mathbf{x}^u)] \right\|, \quad (2)$$

where $D_k(\cdot)$ is the intermediate representation of the discriminator network after the k^{th} layer. Both ground-truth and predicted segmentation masks are concatenated with their corresponding input images. Intuitively, it encourages the generator to predict segmentation maps which have the same feature statistics as the ground truth, and therefore also qualitatively resemble the ground truth. This loss is used on the unlabeled samples \mathbf{x}^u , thus forcing plausible solutions even for cases where dense labels are unavailable.

Self-training loss. During GAN training, the discriminator (D) and the generator (G) networks need to be balanced. If D starts off being too strong, it does not provide any useful learning signal for G . In order to facilitate such balanced dynamics, we introduce the self-training (ST) loss. The main idea is to pick the best generator outputs (i.e. those able to fool D) which do not have the corresponding ground truth, and reuse them for supervised training. Intuitively, this pushes G more to produce predictions which D cannot

distinguish from the real ones. This impedes the progress of D and does not allow it to become too strong.

Technically, the output of D varies between 0 and 1, where 0 should be assigned to the predicted segmentation maps and 1 to the ground-truth segmentation maps. We use this score as a confidence measure for the quality of the predicted segmentations. High-quality predictions are used for supervised training, i.e. we calculate the standard cross-entropy loss based on them. The self-training loss term L_{st} is thus defined as:

$$L_{st} = \begin{cases} - \sum_{h,w,c} \mathbf{y}^* \log S(\mathbf{x}^u), & \text{if } D(S(\mathbf{x}^u)) \geq \gamma \\ 0, & \text{otherwise,} \end{cases} \quad (3)$$

where γ is the confidence threshold which controls how certain D needs to be about the prediction in order for it to be used in self-training; \mathbf{y}^* are the pseudo pixel-wise labels generated from the prediction $S(\mathbf{x}^u)$ of the segmentation network.

The final training objective L_S is composed of the three described terms:

$$L_S = L_{ce} + \lambda_{fm} L_{fm} + \lambda_{st} L_{st}, \quad (4)$$

where $\lambda_{fm}, \lambda_{st} > 0$ are the corresponding weights.

3.1.2 Training D

The discriminator network is trained with the original GAN objective as proposed by Goodfellow *et al.* [10]

$$L_D = \mathbb{E}_{(\mathbf{x}^\ell, \mathbf{y}^\ell) \sim \mathcal{D}^\ell} [\log D(\mathbf{y}^\ell \oplus \mathbf{x}^\ell)] + \mathbb{E}_{\mathbf{x}^u \sim \mathcal{D}^u} [\log(1 - D(S(\mathbf{x}^u) \oplus \mathbf{x}^u))], \quad (5)$$

where \oplus denotes concatenation along the channel dimension. Following the original GAN idea, D learns to differentiate between the real \mathbf{y}^ℓ and the fake segmentation masks $S(\mathbf{x}^u)$ concatenated with the corresponding input images.

3.2 Multi-label Semi-supervised Classification

We extend an ensemble-based semi-supervised classification method (Mean Teacher) [37] for semi-supervised multi-label image classification. This model consists of two networks: a student network G and a teacher network H . Both networks receive the same images under different small perturbations. The weights (θ') of the teacher network are the exponential moving average (online ensemble) of the student network's weights (θ). The predictions made by the student model are encouraged to be consistent with the predictions of the teacher model using the consistency loss which is the mean-squared error between the two predictions.

We optimize the student network using the categorical cross-entropy loss L_{ce} for labeled samples \mathbf{x}^ℓ , and using the consistency loss L_{cons} for all available samples ($\mathbf{x}^{u,\ell}$):

$$L_{MT} = \underbrace{- \sum_c \mathbf{z}^\ell(c) \log(G_\theta(\mathbf{x}^\ell)(c))}_{L_{ce}} + \lambda_{cons} \underbrace{\|G_\theta(\mathbf{x}^{(u,\ell)}) - H_{\theta'}(\mathbf{x}'^{(u,\ell)})\|^2}_{L_{cons}}, \quad (6)$$

where \mathbf{x} and \mathbf{x}' are differently augmented images for student and teacher network respectively, \mathbf{z}^ℓ is the multi-hot vector for ground-truth class labels. The parameter $\lambda_{cons} > 0$ controls the weight of the consistency loss in L_{MT} .

3.3 Network Fusion

The two described branches are trained separately. For evaluation, the output of the classification branch simply deactivates the segmentation maps of those classes not present in the input image:

$$S(\mathbf{x})_c = \begin{cases} 0 & \text{if } G(\mathbf{x}_c) \leq \tau \\ S(\mathbf{x})_c & \text{otherwise} \end{cases} \quad (7)$$

where $S(\mathbf{x})_c$ is the segmentation map for class c , $G(\mathbf{x})_c$ is the soft output of the MLMT-branch, and $\tau = 0.2$ is a threshold on that soft output obtained by cross-validation.

4 EXPERIMENTS

The proposed approach was evaluated on the PASCAL VOC 2012 segmentation benchmark, the PASCAL-Context dataset, and the Cityscapes dataset.

4.1 Setup

4.1.1 Datasets

PASCAL VOC 2012. The dataset consists of 20 foreground object classes and one background class. We use the augmented annotation set which consists of 10582 training images and 1449 validation images. The training set contains 1464 images from the original PASCAL data and 9118 extra images from the Segmentation Boundary Dataset (SBD) [11]. The training data augmentations include random resizing, cropping to 321×321 , and horizontal flipping. All the results for the PASCAL VOC dataset are shown on the validation set.

PASCAL-Context. This is a whole scene parsing dataset containing 4,998 training and 5,105 testing images with dense semantic labels. Following the previous work [4], [21], [40], we used semantic labels for 60 most frequent classes including the background class. The training data augmentations were the same as for the PASCAL VOC dataset.

Cityscapes. This is a driving scene dataset with 2975, 500, 1525 densely annotated images for training, validation, and testing, and contains 19 classes. We downsample the original 1024×2048 images by a factor 2. The training data is augmented with random crops of size 256×512 and horizontal flipping. All the results on the Cityscapes dataset are shown on the validation set.

Evaluation Metric. We report mean Intersection-over-Union (mIoU) for all our experiments.

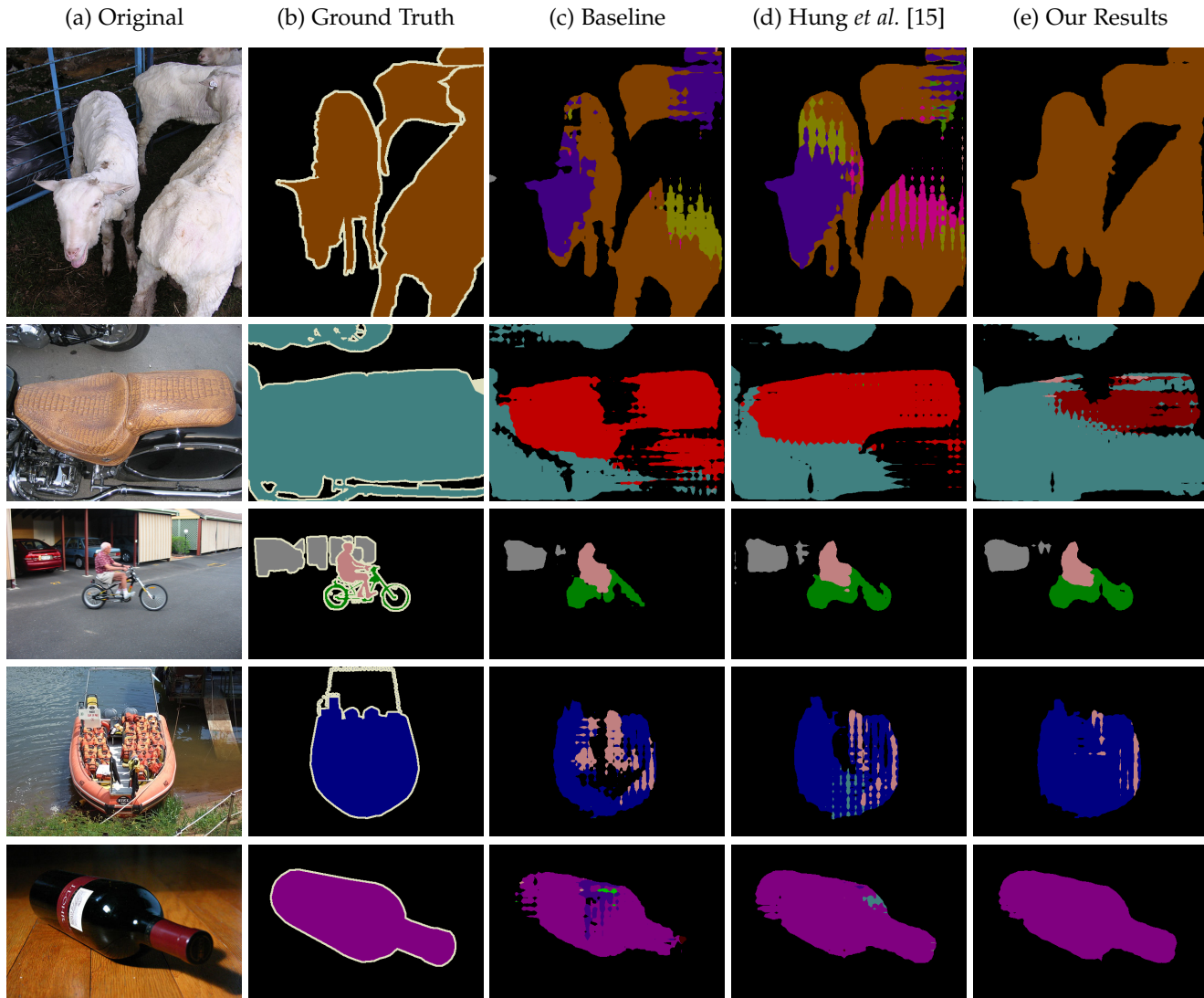


Fig. 4. Qualitative results obtained using our semi-supervised segmentation approach on the PASCAL VOC dataset with 5% labeled data without COCO pre-training.

4.1.2 Network Architecture

Semi-supervised Segmentation GAN. We used DeepLabv2 [4] as our main segmentation network. Due to memory constraints, we used a single-scale variant of it. The discriminator network of the GAN model was a standard binary classification network consisting of 4 convolutional layers with 4×4 kernels with $\{64, 128, 256, 512\}$ channels, each followed by a Leaky-ReLU [25] activation with negative slope of 0.2 and a dropout [35] layer with dropout probability of 0.5. We found this high dropout rate to be crucial for stable GAN training. The last convolutional layer is followed by global average pooling and a fully-connected layer. The output vector representation produced after global average pooling is used for evaluating the feature matching loss.

Semi-supervised Multi-label Classification Network. We used ResNet-101 [12] pre-trained on the ImageNet dataset [8] as the base architecture.

4.1.3 Training details

Similar to [4], we used the poly-learning policy for both the segmentation and the discriminator networks of the

GAN model, where the base learning rate was multiplied by a factor of $((1 - \frac{\text{iter}}{\text{max_iter}})^{\text{pow}})$ in every iteration. In our setup, $\text{pow} = 0.9$. Following the learning scheme in [15], the segmentation network was optimized using the SGD optimizer with a base learning rate of $2.5e-4$, momentum 0.9 and a weight decay of $5e-4$. The discriminator network was optimized using the Adam optimizer [17] with a base learning rate of $1e-4$ and betas set to $(0.9, 0.99)$. The model was trained for 35K iterations on the PASCAL VOC and Cityscapes dataset, and for 50K iterations on the PASCAL-Context dataset. All the learning hyper-parameters remained the same for all datasets except for the Cityscapes dataset, where the base learning rate of the discriminator network was set to $1e-5$. We used a batch size of 8 for both PASCAL datasets and a batch size of 5 for the Cityscapes dataset. Through cross-validation, we find the optimal loss weights: $\lambda_{fm} = 0.1$, $\lambda_{st} = 1.0$, $\lambda_{cons} = 1.0$ and $\tau = 0.2$. These hyper-parameters remained the same for all datasets, whereas we set $\gamma = 0.6$ for both PASCAL datasets and 0.7 for the Cityscapes dataset. Our implementation is based on the open source toolbox Pytorch [30]. All the experiments were run on a Nvidia Tesla P100 GPU.

4.1.4 Baselines

We compare to the DeepLabv2 [4] network as the fully-supervised baseline approach, which was trained only on the labeled part of the dataset. DeepLabv2 makes use of dilated convolutions to enlarge the receptive field size and incorporate larger context, and introduces atrous spatial pyramidal pooling to capture image context at multiple levels.

Our main semi-supervised baseline is the approach proposed by Hung *et al.* [15].

Apart from the differences described in Sec. 2, they also use a two-stage GAN training. In the first stage, both D and G are trained only using labeled data. In the second stage, D 's outputs are used to update G using unlabeled samples, while D itself is further trained only on the labeled images.

4.2 Results

4.2.1 Semi-supervised Semantic Segmentation

We evaluated our approach with different ratios of labeled and unlabeled samples. $1/50$, $1/20$, $1/8$, $1/4$ are the fractions of the total training images in the dataset that are used as labeled data, the rest of the data was used without labels. The labeled samples in the data splits were randomly sampled from the whole dataset, and the same data splits were used for all the baselines.

TABLE 1
Semi-supervised semantic segmentation results on the PASCAL VOC dataset without and with COCO pre-training.

without COCO pre-training				
Method	Labeled Data			
	1/50	1/20	1/8	Full
DeepLabv2	48.3	56.8	62.0	70.7
Hung <i>et al.</i> [15]	49.2	59.1	64.3	71.4
Ours (s4GAN only)	58.1	60.9	65.4	71.2
Ours (s4GAN + MLMT)	60.4	62.9	67.3	73.2
with COCO pre-training				
DeepLabv2	53.2	58.7	65.2	73.6
Hung <i>et al.</i> [15]	57.2	64.7	69.5	74.9
Ours (s4GAN only)	60.9	66.4	69.8	73.9
Ours (s4GAN + MLMT)	63.3	67.2	71.4	75.6

PASCAL VOC Dataset. Table 1 shows the segmentation results on the PASCAL VOC dataset with and without pre-training on the Microsoft COCO [22] dataset. We achieve improved results compared to the previous method for all data splits. Our method achieves a performance increase of 5% to 12% over the baseline for different data splits by utilizing unlabeled samples without pre-training the network on any segmentation dataset. Notably, the approach works well even with only 2% ($1/50$) of labeled data. Figure 4 shows qualitatively how our method helps remove artifacts produced by other methods. We also validated our approach with COCO pre-training to directly compare with Hung *et al.* [15], and achieved an improvement of 6.1 mIoU points over them for the $1/50$ split. We speculate that [15] is inferior in the low-data regime due to the two-stage GAN training, where the discriminator is only updated based on the labeled samples. This effectively reduces the amount



Fig. 5. Qualitative results on the PASCAL-Context dataset using $1/8$ labeled samples. Our approach produces improved results compared to the baseline. We compare our ('Ours') results with the fully-supervised baseline which is trained only on the labeled subset of data.

of data it sees during training, which can easily lead to overfitting.

We conducted our initial experiments using Deeplabv3+ as the backbone architecture. Deeplabv3+ is unstable in the low-data ‘supervised only’ setting. It is only superior, if there is much labeled data. Thus, for a more informative experiment, we rather used Deeplabv2. However, our semi-supervised model achieves even better performance with Deeplabv3+ than with the DeepLabv2-based model, see Table 2.

TABLE 2
Results on PASCAL VOC without COCO pre-training using different backbone architectures.

Method	1/50	1/20	1/8	Full
Deeplabv2 (v2)	48.3	56.8	62.0	70.7
Ours v2 (s4GAN+ MLMT)	60.4	62.9	67.3	73.2
Deeplabv3+ (v3+)	unstable	unstable	63.5	74.6
Ours v3+ (s4GAN+ MLMT)	62.6	66.6	70.4	74.7

The results were obtained with cross-validation to avoid hyper-parameter search on the evaluation set. We also submitted our results to the PASCAL test server. Due to the benchmark restrictions we could only submit one random split (5% labeled samples). The results are consistent with our previous conclusions: 50.1 mIoU for baseline DeepLabv2 vs 60.5 for our semi-supervised method.

PASCAL-Context Dataset. Our approach successfully generalizes to the whole scene parsing PASCAL-Context dataset. Table 3 shows the performance on two splits (1/8 and 1/4 labeled data) of PASCAL-Context. Although this dataset is smaller and more difficult than PASCAL VOC, there is still an improvement over the baseline of 3.2% and 2.4% for the 1/8 and 1/4 splits, respectively.

Fig. 5 show qualitative results on the PASCAL-Context test set using 1/8 labeled samples and the remaining unlabeled samples. PASCAL-Context is a smaller and harder dataset as compared to PASCAL VOC, therefore the results are not as visually appealing. Still, there is a clear improvement over the baseline.

TABLE 3
Semi-supervised semantic segmentation results on the PASCAL-Context dataset without COCO pre-training.

Method	Labeled Data		
	1/8	1/4	Full
DeepLabv2	32.1	35.4	41.0
Hung <i>et al.</i> [15]	32.8	34.8	39.1
Ours (s4GAN only)	34.4	37.1	40.8
Ours (s4GAN + MLMT)	35.3	37.8	41.1

Cityscapes Dataset. On the Cityscapes dataset, the s4GAN branch yields an improvement over the baseline of 3.1% and 1.7% for the 1/8 and 1/4 data splits respectively; see Table 4. The distribution of different classes in this dataset is highly imbalanced. The vast majority of the classes are present in almost every image, and the few remaining classes occur only scarcely. In this situation, a classifier that eliminates labels of non-existing classes does not help, thus, our MLMT branch was ineffective for the Cityscapes dataset.

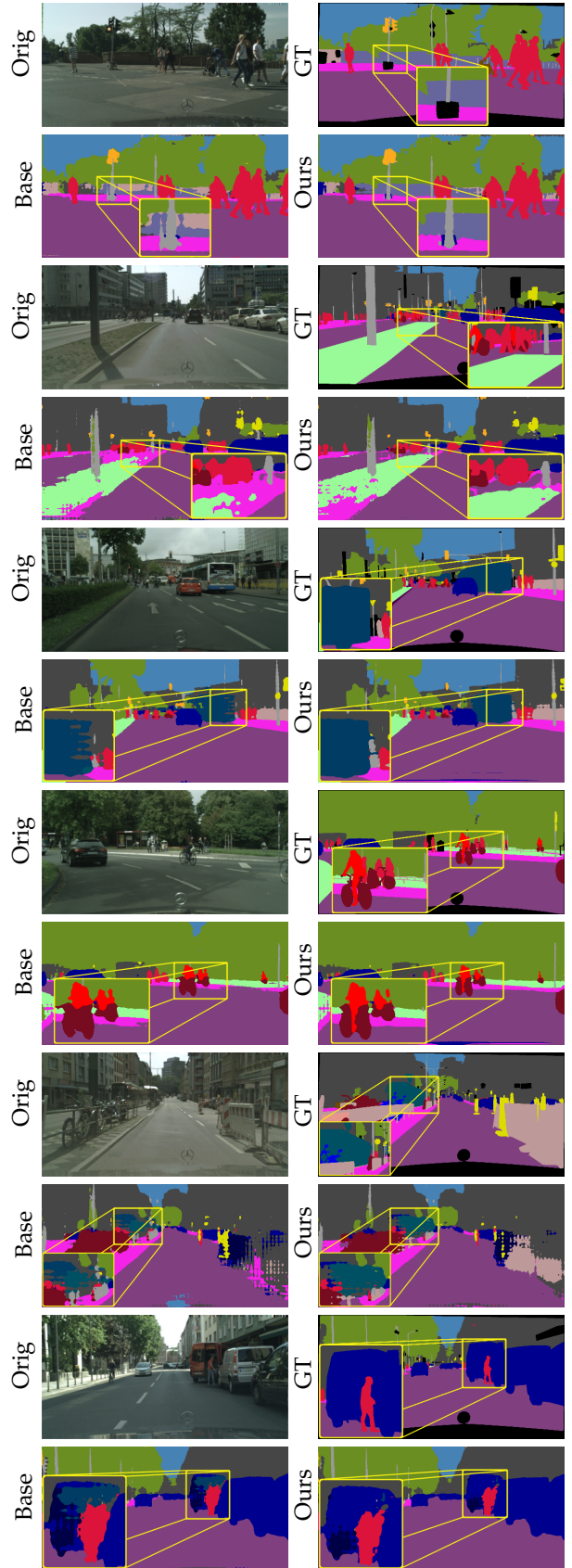


Fig. 6. Qualitative results on the Cityscapes dataset using 1/8 labeled samples without COCO pre-training. The proposed semi-supervised approach produces improved results compared to the baseline. We compare our (‘Ours’) results with the fully-supervised baseline (‘Base’) which is trained only on the labeled subset of data.

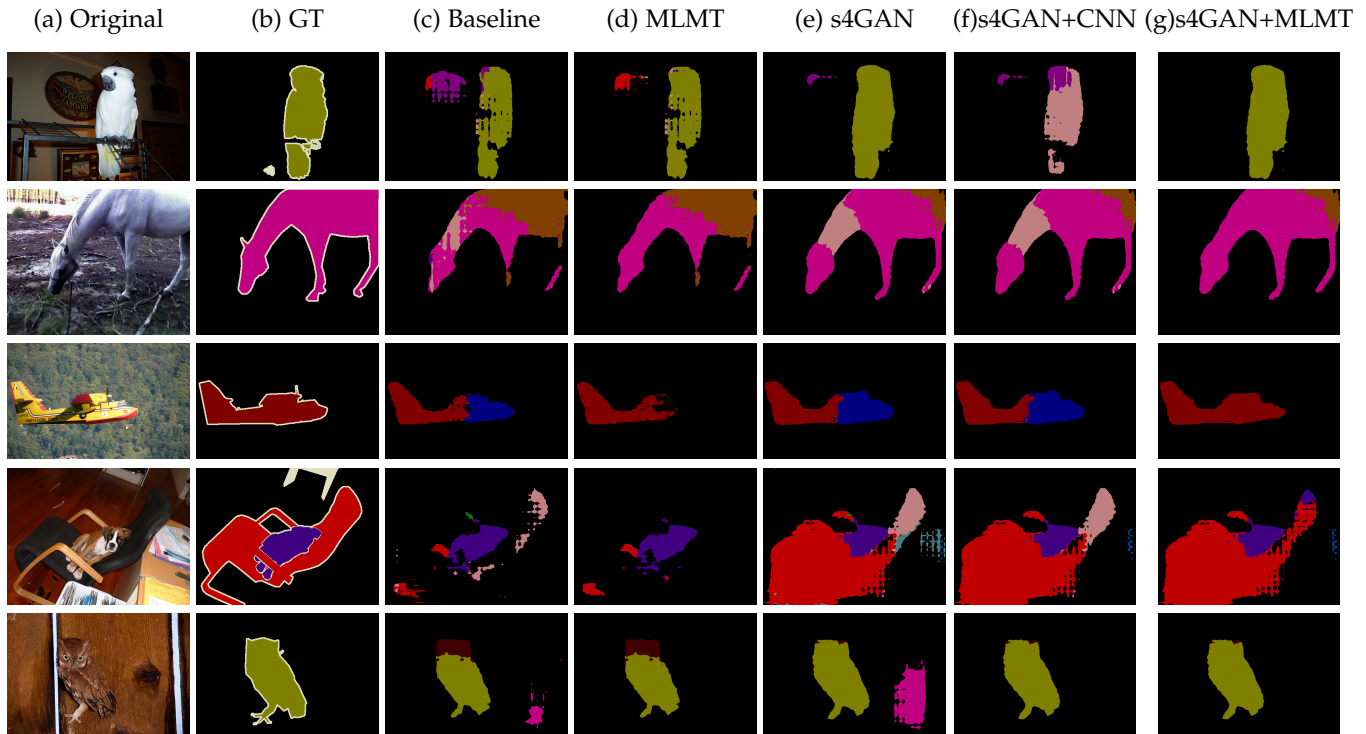


Fig. 7. Ablation study on the PASCAL VOC dataset showing the contribution of the MLMT (d) and the s4GAN (e) branches individually. The s4GAN and the MLMT branches together show a complementary behaviour fixing both low and high-level artifacts (g). These results are obtained using 5% labeled data.

TABLE 4
Semi-supervised semantic segmentation results on the Cityscapes dataset without COCO pre-training.

Method	Labeled Data		
	1/8	1/4	Full
DeepLabv2	56.2	60.2	66.0
Hung <i>et al.</i> [15]	57.1	60.5	66.2
Ours (s4GAN only)	59.3	61.9	65.8

Fig. 6 show qualitative results obtained using our approach with 1/8 labeled samples and the remaining unlabeled samples. The differences on the Cityscapes dataset are subtle, therefore we include the zoomed-in views of informative areas. On images from Fig. 6 show our approach yields improvement over the baseline.

4.2.2 Ablation Studies

All the experiments for the ablation studies are shown on the PASCAL VOC dataset without COCO pre-training.

Contribution of the two branches.

Table 5 shows the contribution of the s4GAN branch and the MLMT branch. The s4GAN branch is able to extract extra dense information using unlabeled images. It improves the shape of the segmented objects, makes the segmentation prediction more coherent by filling small holes, and improves the fine boundaries between the foreground and background. We showcase these improvements in Figure 7(e).

The MLMT branch plays a complementary role and removes the false positives from the predictions. Figure 7(d)

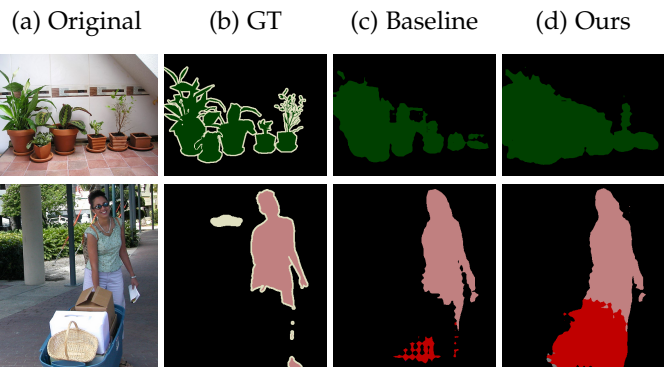


Fig. 8. Failure cases. Sometimes, our approach can lead to under-segmentation of objects with multiple protrusions, shown in the top row. Bottom row shows a case where an ambiguous foreground object is falsely marked as one of the classes.

shows the improvement using the ‘MLMT branch only’ with the segmentation baseline method and Figure 7(g) shows the improvement using the MLMT branch together with the s4GAN branch. The MLMT branch makes use of unlabeled images to extract image-level information about the presence of the certain classes in the image. For some cases, the s4GAN branch introduces new artifacts which are also filtered out by the MLMT branch. This effect is shown in the bottom-row example of Figure 7.

In certain situations our method produces imprecise predictions. Sometimes object classes with multiple protrusions like plant leaves, chair legs, etc. are under-segmented by the s4GAN branch, as shown in Figure 8(top). Occasionally, our

approach can identify certain ambiguous foreground objects as one of the classes, as shown in Figure 8(bottom). Also, there exist few cases where some true positive results are wrongly predicted by the classifier. However, both qualitative and quantitative results confirm that these failure cases are outweighed by the positive effect of the proposed techniques. In Fig. 9, we include a few failure cases for PASCAL-context dataset using our approach. Fig.10 shows a few failure cases for Cityscapes dataset where few thin objects were not segmented properly using our approach.

TABLE 5

Ablation study of the contribution of each branch. Results are shown for the 5:95 data split on the PASCAL VOC dataset.

Method	Data		mIoU
	labeled(%)	unlabeled(%)	
DeepLabv2	5	None	56.8
s4GAN only	5	95	60.9
MLMT only	5	95	59.0
s4GAN + Threshold	5	95	61.2
s4GAN + Class-wise Threshold	5	95	61.5
s4GAN + CNN	5	95	62.2
s4GAN + MLMT	5	95	62.9

Different s4GAN branch loss terms. We trained the generator network with a combination of the cross-entropy (CE) loss, the feature matching (FM) loss, and the self-training (ST) loss.

To justify this configuration, we compare the system performance when using different loss terms; see Table 6. There is a consistent performance increase when adding all the proposed loss terms. We found it crucial for the system stability to train using the FM loss and not the standard GAN loss.

TABLE 6

Ablation study of different GAN loss terms for the generator on the PASCAL VOC dataset. SGAN refers to the standard GAN loss [10], FM refers to the feature-matching loss and ST refers to the self-training loss.

Loss Terms	Labeled Data		
	1/50	1/20	1/8
CE only	48.3	56.8	62.0
CE + SGAN [10]	54.0	57.1	62.5
CE + FM	55.4	58.4	63.9
CE + FM + ST	58.1	60.9	65.4

Figure 11 illustrates the effect of using our proposed self-training loss. We plot how the discriminator score changes during the course of training. The scores are averaged over 100 iterations of fake (generated) and real (ground-truth) samples separately. As discussed in Sec. 3.1.1, adding the ST loss impedes the progress of the discriminator and does not allow it to become overly confident, that is, draws its predicted scores towards 0.5. This has a positive effect on the generator performance, in particular with few labeled samples, as can be seen from the last line of Table 6.

Semi-supervised multi-label classification. In this experiment, we compared the performance of the proposed MLMT branch with a standard supervised classifier. Table 5 shows that we already get an improvement of 1.3%

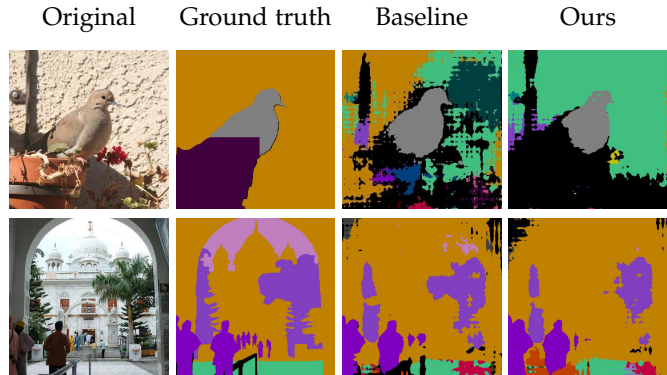


Fig. 9. Qualitative results on the PASCAL-Context dataset using 1/8 labeled samples. Failure of our approach. We compare our ('Ours') results with the fully-supervised baseline which is trained only on the labeled subset of data.

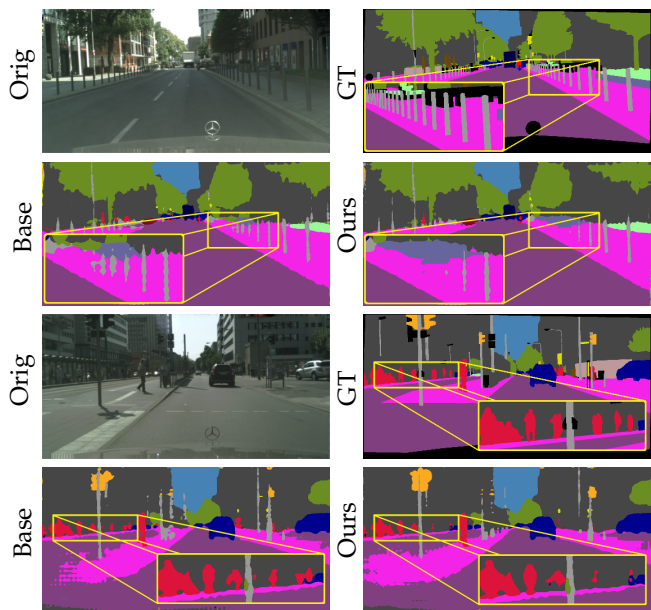


Fig. 10. Qualitative results on the Cityscapes dataset using 1/8 labeled samples. Failures of our approach. We compare our ('Ours') results with the fully-supervised baseline ('Base') which is trained only on the labeled subset of data.

over the s4GAN performance just by using a CNN-based classifier [12], but when we further add the consistency-based semi-supervised classification approach, we observe that the performance improvement increases to 2%. More detailed comparison between the two classification methods is included in the supplementary file. We also conducted a simple heuristic experiment where we deactivate the predicted class channels which have pixel count less than a threshold. In Table 5, 's4GAN + Threshold' refers to the case where a single threshold is set for all the classes and 's4GAN + Class-wise Threshold' refers to the case where each class has its best respective threshold. We search for the best performing thresholds on the validation set in the range from 1K to 12K pixels at an increment step of 1K. Figure 7(f) and (g) show the effect of adding a CNN-based classifier and an MT-based semi-supervised classifier respectively.

We also analyze the performance of the CNN-based

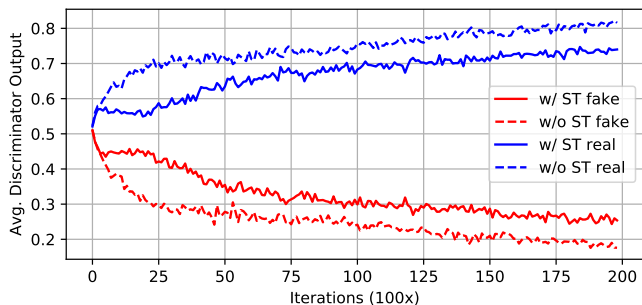


Fig. 11. Evolution of the discriminator output during the course of training averaged over real and fake samples separately. Using the self-training loss (w/ ST) prevents D from becoming overly strong and results in better training dynamics compared to the case when self-training is disabled (w/o ST).

multi-label classification and MLMT-based semi-supervised multi-label classification independent of the segmentation model. Figure 12 shows the comparison between the ROC curves of the two methods on the task of multi-label classification. The MLMT classifier obtains a lower false positive rate for the same true positive rate. The effect is even more pronounced when not using ImageNet pre-training; see Figure 12(b). This mode of operation is important for domains where ImageNet pre-training does not help, e.g. bio-medical image analysis.

TABLE 7

Semi-supervised semantic segmentation results on the PASCAL VOC dataset using extra weak image-level annotations. Data splits ($A/B/C$) refers to the usage of A pixel-wise labeled samples, B image-level labeled samples and C unlabeled samples.

Method	Data Split (Strong/Weak/Unlab)		
	1.4K/0/9K	1.4K/9K/0	All/0/0
DeepLab-CRF-LargeFOV [3]	62.5 ^c	—	67.6 ^c
WSSL (CRF) ^a [29]	—	64.6	—
MDC ^a [39]	—	62.7	—
MDC (CRF) ^a [39]	—	65.7	—
DeepLabv2	65.7	—	70.7
Ours (s4GAN only) ^b	67.5	—	71.2
Ours (s4GAN + MLMT) ^b	69.6	70.9	72.9

^a Base network: DeepLab-LargeFOV,

^b Base network: DeepLabv2, ^c As reported in [29]

4.2.3 Semi-supervised Semantic Segmentation with Weak-labels

To further validate the effectiveness of our approach, we compare it to other semi-supervised segmentation methods [29], [39] that utilize extra weak image-level annotations. Here, we compare the performance of our approach with methods that use extra image-level annotations *i.e.* 1,464 strongly (w/ segmentation masks) annotated images from the original PASCAL VOC dataset and 9,118 weakly (image-level) annotated images from the augmented SBD dataset. To use extra image-level annotations, we train the MLMT branch using extra image-level labels for improved multi-label classification. The training procedure and hyperparameters remain exactly same as in the previous semi-supervised setting. Table 7 summarizes the semi-supervised

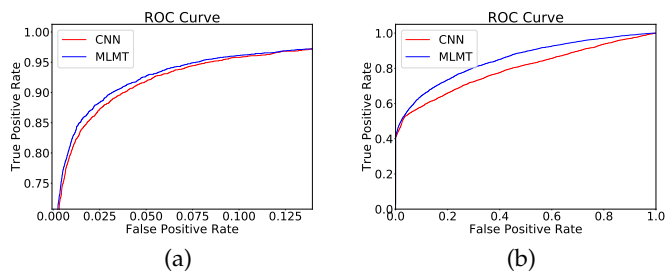


Fig. 12. ROC curves for CNN-based classification and MT-based semi-supervised classification method using 5% labeled data with (a) and without (b) ImageNet pre-training. MT produces fewer false positives, especially when training from scratch.

semantic segmentation results with extra $\sim 9K$ image-level annotations. We achieve an improvement of 5.2% over the baseline. Unlike previous methods, our approach does not utilize the CRF post-processing.

5 CONCLUSION

In this work we presented a two-branch approach to the task of semi-supervised semantic segmentation. The branches are designed to alleviate both low-level and high-level artifacts which often occur when working in a low-data regime. The effectiveness of this design is demonstrated in a series of extensive experiments.

ACKNOWLEDGEMENTS

This study was supported by the German Federal Ministry of Education and Research via the project Deep-PTL and by the Intel Network of Intelligent Systems. We also thank Facebook for their P100 server donation and gift funding.

REFERENCES

- [1] J. Ahn and S. Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *CVPR*, 2018.
- [2] B. Athiwaratkun, M. Finzi, P. Izmailov, and A. G. Wilson. Improving consistency-based semi-supervised learning with weight averaging. *arXiv preprint arXiv:1806.05594*, 2018.
- [3] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *CoRR*, abs/1412.7062, 2014.
- [4] L. 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. *TPAMI*, 40(4):834–848, 2018.
- [5] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018.
- [6] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [7] J. Dai, K. He, and J. Sun. Boxesup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *ICCV*, 2015.
- [8] J. Deng, W. Dong, R. Socher, L. jia Li, K. Li, and L. Fei-fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [9] M. Everingham, L. Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2):303–338, 2010.
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*. 2014.

- [11] B. Hariharan, P. Arbelaez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse detectors. In *ICCV*, 2011.
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [13] S. Hong, H. Noh, and B. Han. Decoupled deep neural network for semi-supervised semantic segmentation. In *NIPS*. 2015.
- [14] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *CVPR*, 2018.
- [15] W. Hung, Y. Tsai, Y. Liou, Y. Lin, and M. Yang. Adversarial learning for semi-supervised semantic segmentation. In *BMVC*, 2018.
- [16] A. Khoreva, R. Benenson, J. Hosang, M. Hein, and B. Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *CVPR*, 2017.
- [17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [18] S. Laine and T. Aila. Temporal ensembling for semi-supervised learning. In *ICLR*, 2017.
- [19] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324, 1998.
- [20] D. Lin, J. Dai, J. Jia, K. He, and J. Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *CVPR*, 2016.
- [21] G. Lin, A. Milan, C. Shen, and I. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *CVPR*, 2017.
- [22] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.
- [23] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- [24] P. Luc, C. Couprie, S. Chintala, and J. Verbeek. Semantic segmentation using adversarial networks. In *NIPS Workshops*. 2016.
- [25] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML Workshop on Deep Learning for Audio, Speech and Language Processing*, 2013.
- [26] T. Miyato, S. ichi Maeda, M. Koyama, and S. Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *TPAMI*, 2018.
- [27] 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.
- [28] A. Oliver, A. Odena, C. Raffel, E. D. Cubuk, and I. J. Goodfellow. Realistic evaluation of semi-supervised learning algorithms. In *ICLR Workshop*. 2018.
- [29] G. Papandreou, L.-C. Chen, K. P. Murphy, and A. L. Yuille. Weakly- and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *ICCV*, 2015.
- [30] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017.
- [31] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun. Large kernel matters – improve semantic segmentation by global convolutional network. In *CVPR*, 2017.
- [32] P. O. Pinheiro and R. Collobert. From image-level to pixel-level labeling with convolutional networks. In *CVPR*, 2015.
- [33] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, X. Chen, and X. Chen. Improved techniques for training gans. In *NIPS*. 2016.
- [34] N. Souly, C. Spampinato, and M. Shah. Semi supervised semantic segmentation using generative adversarial network. In *ICCV*, 2017.
- [35] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *JMLR*, 15(1):1929–1958, 2014.
- [36] M. Tang, F. Perazzi, A. Djelouah, I. B. Ayed, C. Schroers, and Y. Boykov. On regularized losses for weakly-supervised cnn segmentation. In *ECCV*, 2018.
- [37] A. Tarvainen and H. Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NIPS*. 2017.
- [38] Y. Wei, J. Feng, X. Liang, M.-M. Cheng, Y. Zhao, and S. Yan. Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In *CVPR*, 2017.
- [39] Y. Wei, H. Xiao, H. Shi, Z. Jie, J. Feng, and T. S. Huang. Revisiting dilated convolution: A simple approach for weakly- and semi-supervised semantic segmentation. In *CVPR*, 2018.
- [40] H. Zhang, K. Dana, J. Shi, Z. Zhang, X. Wang, A. Tyagi, and A. Agrawal. Context encoding for semantic segmentation. In *CVPR*, 2018.