

A Domain Agnostic Normalization Layer for Unsupervised Adversarial Domain Adaptation

R. Romijnders P. Meletis G. Dubbelman
Eindhoven, University of Technology
romijndersrob@gmail.com

Abstract

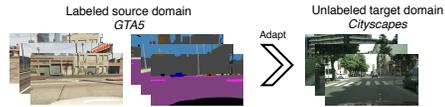
We propose a normalization layer for unsupervised domain adaption in semantic scene segmentation. Normalization layers are known to improve convergence and generalization and are part of many state-of-the-art fully-convolutional neural networks. We show that conventional normalization layers worsen the performance of current Unsupervised Adversarial Domain Adaption (UADA), which is a method to improve network performance on unlabeled datasets and the focus of our research. Therefore, we propose a novel Domain Agnostic Normalization layer and thereby unlock the benefits of normalization layers for unsupervised adversarial domain adaptation. In our evaluation, we adapt from the synthetic GTA5 data set to the real Cityscapes data set, a common benchmark experiment, and surpass the state-of-the-art. As our normalization layer is domain agnostic at test time, we furthermore demonstrate that UADA using Domain Agnostic Normalization improves performance on unseen domains, specifically on Apolloscape and Mapillary.¹

1. Introduction

Semantic segmentation constitutes a crucial task in computer vision of assigning a class to every pixel in an image. Applications range from autonomous driving and robotic navigation to segmenting natural scenes. Convolutional neural networks (CNNs) have shown good performance, e.g. [4, 32]. However, when we train these models on a *source domain*, and evaluate on another, *target domain*, the performance degrades. Ideally, we would adapt our model to the target domain, without requiring to label the images (unsupervised). In this work, we aim for unsupervised domain adaptation to improve generalization capability for semantic scene segmentation.

Unsupervised domain adaptation addresses three problems: 1) In new domains, images differ in appearance, light-

Domain Agnostic training



Domain Agnostic testing

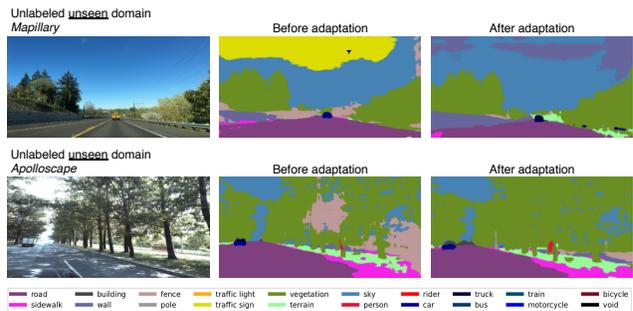


Figure 1. Diagram for unsupervised domain adaptation. During training, we only have labels for the source domain. During testing, we want a model that we can apply to any unseen and unlabeled domain, without any domain-specific fine tuning.

ing, contrast or colorization. For example, when the time of day, season or camera changes. This shift is studied in [23, 39, 38]. We want good performance for all domains; 2) Labeling the data is a cumbersome task. For example, labeling the Cityscapes data set took up to 90 minutes per image [12]. Unsupervised domain adaptation alleviates this labeling burden as no labels are required for the target domain; 3) Models trained on simulated data fail to perform well on real data. Domain adaptation poses a potential solution to close this gap.

A recent and promising work in unsupervised domain adaptation for semantic segmentation is unsupervised adversarial domain adaptation (UADA) [15], which will be the focus of our work. Current approaches in UADA have two challenges. First, evaluations of the model focus on the source and target domain, but omit any unseen domains. For example, when we adapt a segmentation model from Germany to China, we also want good performance in other Asian countries. To summarize, we aim to improve the gen-

¹Accepted to IEEE WACV 2019

eralization capability of the model to any (unseen) domain, not limiting to only the source and target domain.

As a second challenge, many recent approaches in UADA use CNNs without normalization layers. This is sub-optimal as normalization layers speed up convergence and reduce sensitivity to initialization of the parameters and hyperparameters [22, 24, 36, 8]. For example, the winning entries of the recent ImageNet competition and Robust Vision competition at CVPR 2018 use normalization layers [27, 5, 45, 20]. However, none of previous work in UADA use normalization layers [35, 19, 10]. This absence probably follows from observations that normalization layers reduce performance during adversarial training [34, 44].

Batch normalization [22], the most commonly used normalization layer, makes the output for one image dependent on another image in a batch. This dependency creates havoc when the batch contains images from both labeled and unlabeled domains. Therefore, we propose a new normalization layer for UADA. In Section 5.1, we show that our normalization layer does not degrade performance in UADA, as conventional normalization layers do. Next, in Section 5.2, we show that we surpass state-of-the-art, [19], adaptation using our normalization layer. Moreover, we show that this adaptation also improves performance on two unseen domains.

We evaluate our approach on large-scale scientific benchmarks for semantic scene segmentation. We experiment on data from urban scenes as these data pose many challenges due to object clutter and have a wide diversity of classes and appearance. We consider two adaptation tasks: from synthetic to real data, GTA5 to Cityscapes; from real to real data, Cityscapes to Mapillary [12, 33, 30]. To the best of our knowledge, we are the first to evaluate performance on unseen domains, Mapillary and Apolloscape [30, 21].

As we unlock the benefits of normalization layers for UADA, we make the following contributions:

- We demonstrate a degraded performance when using conventional batch normalization in UADA and provide insight in this degradation using an experiment. (Section 5.1)
- We propose a Domain Agnostic Normalization layer for UADA (Section 3.2), surpass state of the art, [19], and show a performance improvement on two unseen domains (Section 5.2).

Our code to reproduce the results is available at [43].

2. Related work

2.1. Semantic segmentation

Semantic segmentation has been widely studied in computer vision. Many current semantic segmentation models

use CNNs, following the progress in large-scale image classification [14, 17]. Neural networks learn hierarchical representations using layers of neurons [7].

Recently, fully convolutional networks (FCN) generalize CNNs for arbitrary input sizes [29]. This generalization enables the use of classification networks such as VGG [37] or ResNet [18] for semantic segmentation. In this generalization, each layer comprises a grid of individual representations, one representation for each receptive field. As semantic segmentation deals with large images, [46] introduced dilated convolutions to enlarge the receptive field and the authors of [4, 31] proposed to learn the upsampling for large label maps.

2.2. Adversarial adaptation

Adversarial adaptation draws inspiration from Generative Adversarial Networks (GAN) [16]. In training a GAN, a generator outputs data samples while trying to confuse a discriminator that classifies between generated and real data samples.

Analogous to using a discriminator to align generated and real data samples, adversarial training uses an adversary to align representations from multiple domains in a neural network. In other words, a domain classifier serves as an adversary and learns to discriminate source from target representations. Another model learns to confuse the domain classifier, which encourages alignment between the source and target representations. As the representations from multiple domains become more aligned, the neural network will generalize better to new domains [6].

2.3. Domain adaptation

Ganin et al. [15] propose adversarial adaptation to align the representations. This work has been extended in [40] and [41] for image classification. Another adversarial approach to domain adaptation is to transfer the style of an image from source to target domain [28, 9, 48].

UADA has been applied to semantic segmentation. For example, [19] builds on the work of [40] and applies a domain classifier on the representations learned by an FCN. UADA aims for alignment at the representation level. Later works have focused on combining alignment at the representation level with alignment at the logit level, [10], or alignment at the output level, [35]. We focus on UADA for semantic segmentation models that use normalization layers. Normalization layers pose specific challenges for UADA that we will address in this work.

2.4. Normalization

Using normalization layers in a neural network speeds up training and reduces sensitivity to initialization of the parameters and hyperparameters [22, 24, 36, 8].

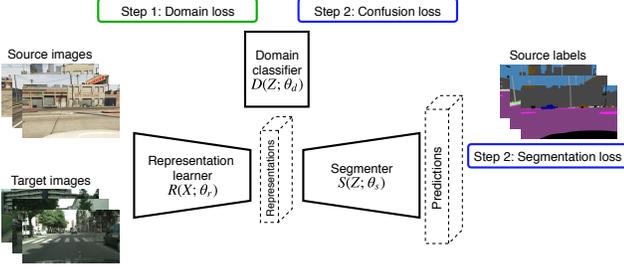


Figure 2. Diagram depicting the two alternating steps during UADA. In the first step, we train the domain classifier using the domain loss. In the second step, we train the segmentation model using the segmentation and confusion loss. (see Equation 2)

Perhaps the earliest attempt was local response normalization by Krizhevsky et al. [25]. Ioffe and Szegedy introduced the batch normalization layer that normalizes all representations in a batch [22]. During training, the internal activations in the many layers may shift. Batch normalization normalizes the inputs to zero mean and unit variance.

Batch normalization has the disadvantage that predictions depend on all samples in a batch. To alleviate this dependency, other normalization layers, such as instance normalization [42] and layer normalization [3] are proposed. For a representation grid of size $N(\text{Sample size}) \times H(\text{Height}) \times W(\text{Width}) \times C(\text{Channels})$ in a CNN, these layers normalize over the following axes:

- *Batch normalization* normalizes over N, H, W
- *Instance normalization* normalizes over H, W
- *Layer normalization* normalizes over H, W, C

In this work, we compare our proposed normalization layer with two forms of batch normalization and with instance normalization. We do not consider layer normalization, as batch normalization is known to outperform layer normalization in CNNs [3]. In Section 3.2, we propose our Domain Agnostic Normalization, compare to other layers in Section 5.1, and show the benefits for domain agnostic testing in Section 5.2.

3. Method: Domain Agnostic Normalization for unsupervised domain adaptation

In this section, we present the details of our proposed Domain Agnostic Normalization layer. First, in Section 3.1, we introduce notation and set up UADA. Second, in Section 3.2, we address the problems of conventional normalization in UADA and propose Domain Agnostic Normalization.

3.1. Domain adaptation

We aim to learn a segmentation model that takes an image, $X \in \mathbb{R}^{H \times W \times C}$, and segments the image to a label

map, $\hat{Y} \in \mathbb{R}^{H \times W \times K}$. Here, C is the number of input channels, K is the number of output classes and H and W are height and width, respectively. We train using labeled source data, (X_s, Y_s) , and unlabeled target data, X_t .

Figure 2 displays a diagram of the UADA model. The model consists of three major blocks. First, the representation learner parametrizes the representations using a CNN, $Z = R(X; \theta_r)$. This representation learner can contain an arbitrary number of normalization layers. We learn the parameters of the representation learner using both the source (labeled) and target (unlabeled) domain. Second, the segmenter maps the representations to a class prediction for each pixel, $\hat{Y} = S(Z; \theta_s)$. We learn the parameters of the segmenter using the source domain. We assume that all domains have the same classes. So during testing, we re-use the segmenter on any domain (Domain Agnostic testing, Figure 1). Finally, the domain classifier assigns a probability, p , that a representation comes from an image in the target domain, $p = D(z; \theta_d)$. We learn the parameters of the domain classifier using the images from both the source and target domain.

Training UADA requires two alternating steps as confusing the representations opposes the learning of a domain classifier. In the first of alternating steps, we minimize the domain loss, $\mathcal{L}_{dom}(X_s, X_t; \theta_r, \theta_d)$, to learn the parameters of the domain classifier. As we have two domains, the domain loss is a binary cross entropy loss. In the second of alternating steps, we minimize the segmentation loss, $\mathcal{L}_{segm}(X_s, Y_s; \theta_r, \theta_s)$, and the confusion loss, $\mathcal{L}_{conf}(X_t; \theta_r, \theta_d)$. The segmentation loss is a multi-class cross entropy for the K output classes. The confusion loss encourages the segmentation model to confuse a target representation for a source representation. In other words, the confusion loss assumes a low value when the domain classifier predicts a low value of p for any target representation. Therefore, the confusion loss follows:

$$\mathcal{L}_{conf}(X_t; \theta_r, \theta_d) = - \sum_{z \in R(X_t; \theta_r)} \log(1 - D(z; \theta_d)) \quad (1)$$

In total, we train unsupervised adversarial domain adaptation using the following alternating steps:

$$\begin{aligned} & \min_{\theta_d} \mathcal{L}_{dom}(X_s, X_t; \theta_r, \theta_d) \\ & \min_{\theta_r, \theta_s} \mathcal{L}_{segm}(X_s, Y_s; \theta_r, \theta_s) + \lambda \mathcal{L}_{conf}(X_t; \theta_r, \theta_d) \end{aligned} \quad (2)$$

λ trades off the segmentation loss and the confusion loss. Figure 2 displays a diagram for this alternating scheme.

3.2. Domain Agnostic Normalization layer

In this section, we propose our Domain Agnostic Normalization (DAN) layer for UADA. First, we outline a major problem with batch normalization, the most commonly

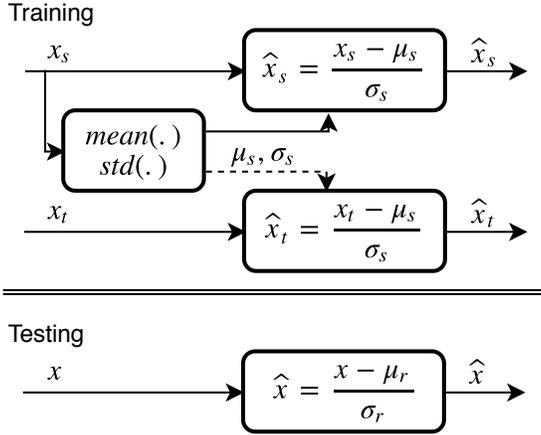


Figure 3. Diagram for our Domain Agnostic Normalization layer. We transform the source representations, x_s , using the source statistics, μ_s and σ_s . We transform the target domain and all unseen domains using the source statistics as fixed parameters. Hence, the dotted line indicates that no gradients flow during back propagation.

used normalization layer [22]. Second, we address these problems and propose DAN.

Batch normalization introduces dependencies between all representations in a batch. For each batch, batch normalization calculates the statistics on all representations and applies these statistics to any representation. This dependency makes the supervision on any image depend on the representations of any other image in the batch. When doing multi-domain training, batches contain images of multiple domains. Now this dependency makes the supervision on one domain depend on representations in another domain. For multi-domain training using only labeled domains, this dependency poses no problem. For a representation in one image, the supervision from its own prediction would be stronger than the implicit supervision of other images via the statistics. For multi-domain training also using an unlabeled domain, as in UADA, however, this dependency creates havoc on the representations of the unlabeled domain. These representations get only supervision from other images via the statistics and have no supervision from a prediction to outweigh this influence. Such a learning procedure would counteract segmentation performance on the target domain. Our experiments in Section 4 will confirm this detrimental effect.

The incompatibility of conventional batch normalization and UADA motivates our requirements for a new normalization layer:

1. A normalization layer that introduces no dependency between source and target representations. Learning representations can only depend on information from

one domain.

2. A normalization layer that applies the same transformation to any domain. Minimizing the confusion loss in Equation 1 should concern the differences in learned representations and not the differences in normalization transformations. Moreover, as we use the same transformation on any domain, we do not need to fine-tune our model when we test on unseen domains (Domain Agnostic testing, Figure 1).

To satisfy these requirements, we propose our Domain Agnostic Normalization layer. Figure 3 shows a diagram of DAN. During training, DAN makes the following steps:

1. For the *source domain*, our normalization layer operates like batch normalization. We calculate the statistics on all representations of the source domain per batch and transform the representations using these statistics.
2. For the *target domain*, our normalization layer transforms the representations using the statistics from the source as fixed parameters. This fixation of parameters causes no gradients to flow during backpropagation, so our normalization layer introduces no dependencies across domains.

During testing, we use the same transformation on all domains, the source domain, the target domain, and any unseen domains. This transformation uses the statistics for the source domain that we aggregate during training. Other works have proposed to re-estimate the statistics for each new domain [26]. Our DAN remains agnostic to the domain at test time and does not require the additional effort of re-estimating the statistics.

4. Experiments

In this section, we outline three sets of experiments that demonstrate the benefit of DAN in UADA.

Our first set of experiments compares the alternatives for DAN in UADA. We run experiments using four normalization layers:

- *Conventional batch normalization* introduces a dependency between representations from source and target domain. When the target domain has no labels, we expect that batch normalization decreases the performance on the target domain.
- *Split batch normalization* [11] applies one conventional batch normalization layer per domain. This split removes the dependencies between domains. However, split batch normalization does not apply the same transformation to each domain.

- *Instance normalization* [42] normalizes the representations per image and introduces no dependencies at all. Instance normalization, however, calculates the statistics using $H \times W$ representations per layer instead of the $N \times H \times W$ representations that batch normalization uses. This smaller sample introduces more stochasticity.
- *Domain Agnostic Normalization* introduces no dependencies between representations across domains and applies the same transformation in each domain.

Our second set of experiments demonstrates the benefits of DAN on two adaptation tasks:

- Adaptation from synthetic to real data, from GTA5 to Cityscapes. We run this experiment to confirm that we reproduce state of the art in UADA, [19, 10, 35]. We compare to these works in Section 5.1.
- Adaptation from real to real data, from Cityscapes to Mapillary. So far as we know, we are the first work to report on this adaptation task. We value this adaptation task as Cityscapes is shot in one country using one camera, whereas Mapillary is shot in multiple countries using multiple cameras (Section 4.1 provides the details per data set).

Our third set of experiments evaluates the final models of both adaptation tasks on unseen domains. The main benefit of DAN occurs in domain agnostic testing. DAN acts as a fixed transformation during testing. One might wonder if this fixation generalizes to unseen domains. Therefore, we evaluate the segmentation model that we adapt from GTA5 to Cityscapes on two unseen domains, Mapillary and Apolloscape. We evaluate the segmentation model that we adapt from Cityscapes to Mapillary on the unseen domain, Apolloscape. To the best of our knowledge, we are the first to evaluate these adaptation tasks on unseen domains.

4.1. Data sets

Cityscapes[12] contains images from 50 cities in and around Germany. Images show scenery over several months, all in good weather conditions. The dense labels consist of nineteen categories of the urban scene, *e.g.* road, sidewalk, car, bus and traffic sign. All images are shot using the same car and camera. Sample sizes: training, 2975; validation, 500.

GTA5[33] contains images rendered from the Grand Theft Auto computer game. Images are taken in the fictitious city Los Santos. The dense labels follow the class definitions of the Cityscapes data set. Sample sizes: training, 9000; validation, 3000.

Mapillary[30] contains images from all around the world. Images are taken in varying weather conditions,

seasons, and time of day. Images are shot using different devices such as mobile phones, action cameras and professional equipment. The dense labels cover 66 classes, but we use only the 19 classes that overlap with Cityscapes and GTA5. Sample sizes: training, 18000; validation, 2000.

Apolloscape[21] contains images from Beijing, China. Images are taken in bright weather conditions and are shot using the same car and camera. We use only the 19 classes that overlap with Cityscapes and GTA5. The labeling process was partially automated, which introduces artefacts. Sample size: we use 8327 images from the validation set.

4.2. Performance metric

We follow the evaluation in [19, 35, 47] and report results on 19 classes. We consider mean intersection over union (mIOU) as our figure of merit. Intersection over union (IOU) represents the number of pixels predicted correctly per class divided by the total number of predicted and true pixels for that class. mIOU averages the IOU over all classes. Perfect segmentation would achieve 100 % mIOU.

We assess the confusion of source and target representations using a retrieval curve, as in [35]. Per target representation, we consider a number of nearest representations (horizontal axis) and average the number of source representations retrieved (vertical axis). Perfect alignment occurs when half of the nearest representations are source representations and half of the nearest representations are target representations. In other words, for any m target representations, perfect alignment occurs when $\frac{m}{2}$ representations come from the source domain and $\frac{m}{2}$ representations come from the target domain. The more two sets of representations are aligned, the closer their retrieval curve lies to the perfect alignment curve. In all retrieval curves in this work, we sample five thousand source and five thousand target representations from images in the respective validation sets.

4.3. Implementation details

In all models, the segmentation model trains for 17 epochs. Adversarial adaptation learns from two data sets at the same time. We define an epoch when, in expectation, we sampled the smallest data set once. Due to computational limits, we resize all images and labels to size 384×768 . For the purposes of this work, we consider a new data set to be a new domain.

For training, we use stochastic gradient descent with momentum. Learning rate and momentum are set at 0.01 and 0.9, respectively. We half the learning rate after 9 and 16 epochs. The batches contain two images per domain. The network uses pre-trained parameters from ResNet-50 [18], trained on ImageNet [13]. To stabilize the training of the domain classifier, we use instance noise [2]. All experiments run in Tensorflow 1.6 [1].

Table 1. **Analyzing batch normalization on multi-domain training.** The units are % mIOU. The *batch constituents* column indicates the domains of images in the batch, G refers to GTA5 and C refers to Cityscapes. In this experiment, we train models using only the GTA5 labels.

Normalization layer	Batch constituents	Tested on	
		GTA5	Cityscapes
Batch norm	(G, G)	62.1	34.8
Batch norm	(G, G, C, C)	61.2	22.1
DAN	(G, G, C, C)	61.9	35.0

Table 2. **Comparing normalization layers on adaptation GTA5 → Cityscapes.** Numbers represent performance on the target domain. Gap reduction indicates the difference between source-only and UADA. Norm. abbreviation normalization. All units are % mIOU.

Normalization layer	Source only	UADA	Gap reduction
No norm.	26.5	28.4	1.9
Batch norm. [22]	34.8	29.6	-5.2
Instance Norm. [42]	30.3	31.4	1.1
Split batch norm. [11]	34.8	35.4	0.6
DAN [ours]	34.8	38.2	3.4

5. Results

5.1. Comparing normalization layers

First, we report experimental results on single domain training using conventional batch normalization and our normalization layer. In Section 3.2, we reasoned that conventional batch normalization degrades the performance on the unlabeled target domain in unsupervised domain adaptation. This degradation is undesirable for UADA, where we aim for segmentation performance on the target domain.

Table 1 reports the results for this experiment. A source-only model on GTA5 using conventional batch normalization achieves 62.1 % mIOU on the GTA5 and 34.8 % mIOU on Cityscapes. When we include the unlabeled images from Cityscapes in the batches, the performance changes to 61.2 % mIOU on GTA5 and 22.1 % mIOU on Cityscapes. These results confirm that conventional batch normalization degrades the performance on the unlabeled target domain, Cityscapes, by 12.7 points mIOU. When we change the conventional batch normalization layer to our DAN, the model achieves 61.9 % mIOU on GTA5 and 35.0 % mIOU on Cityscapes. These numbers are comparable to single domain training and confirm that multi-domain training using our DAN does not degrade the performance on the unlabeled target domain.

The retrieval curve in Figure 4 shows another view on the effects of batch normalization in unsupervised domain adaptation. We compare the retrieval curves for using conventional batch normalization or DAN when we include im-

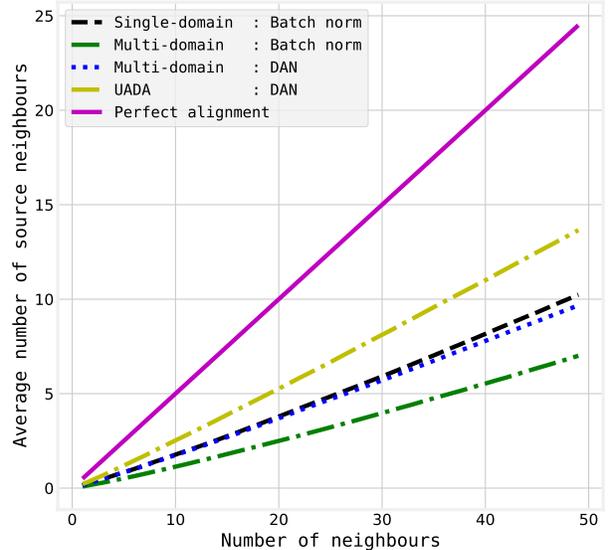


Figure 4. Retrieval curve to show the effect of conventional batch normalization on the representations. Per target representation, we consider a number of nearest representations (horizontal axis) and average the number of source representations retrieved (vertical axis). Perfect alignment occurs when half of the nearest representations are source representations and half of the nearest representations are target representations; the more aligned the representations are, the closer their curve lies to the perfect alignment curve. We explain retrieval curves in Section 4.2

ages from an unlabeled domain in the batch. The curve for conventional batch normalization in multi-domain training (green) lies below the source-only curve (black) and the curve for multi-domain training using DAN (blue). This difference shows that conventional batch normalization learns representations for the unlabeled domain that are less aligned with the representations from the labeled domain. This decrease in alignment opposes unsupervised domain adaptation, where we aim for perfect alignment of the representations in both domains.

Next, we compare DAN with four alternatives: conventional batch normalization, no normalization, split batch normalization, [11], and instance normalization, [42]. Table 2 shows the results for UADA using the five normalization layers. Conventional batch normalization degrades the performance on the target domain, as we also observed in Table 1. Instance normalization introduces more noise during training and we observe that even the performance on source-only is lowest of the three normalization layers. Split batch normalization does not apply the same transformation to each domain and has the lowest gap reduction. Finally, we observe that the gap reduction using DAN ranks highest at 3.4 points mIOU.

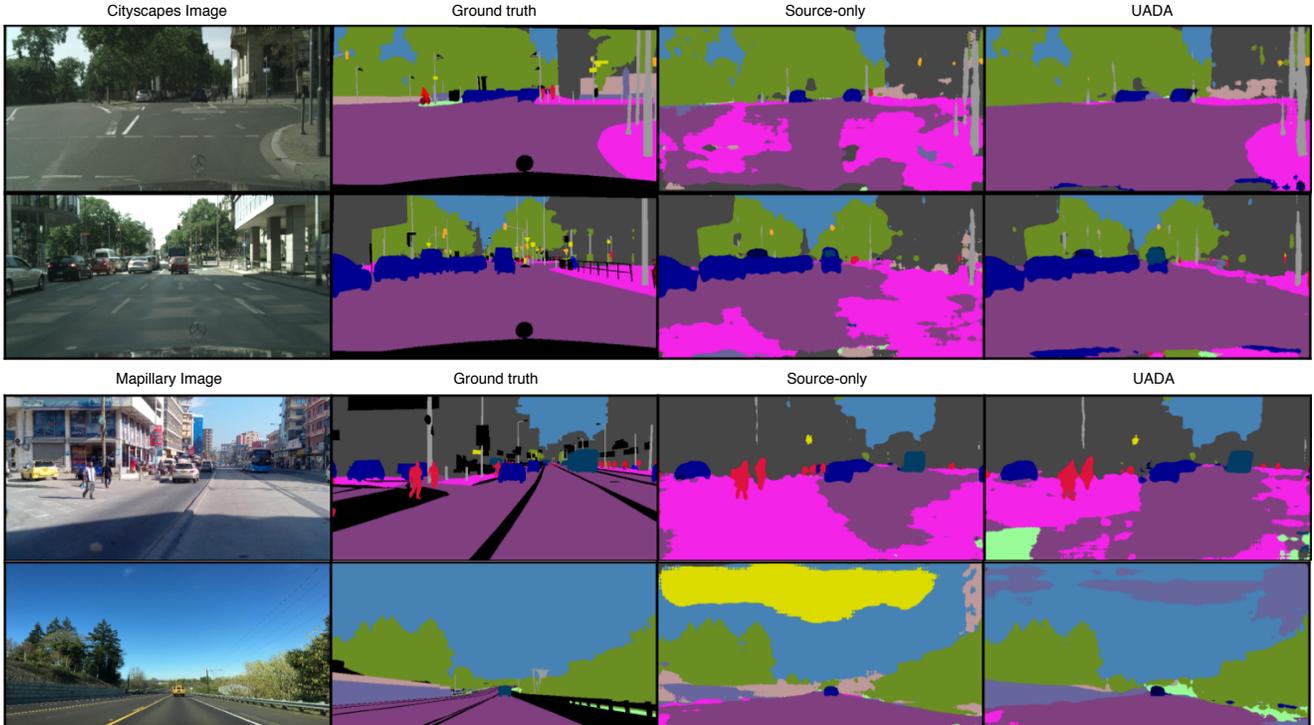


Figure 5. Example results for images from the Cityscapes and Mapillary validation sets. Per image, we show predictions from a model trained on source-only (GTA5) and UADA ($GTA5 \rightarrow Cityscapes$) using DAN. The supplementary material contains more example results.

Table 3. **Hyperparameter analysis** Comparing target performance for ranging values of λ . CS abbreviates Cityscapes. Numbers represent target performance in units % mIOU

value of λ	10^{-5}	10^{-4}	10^{-3}	10^{-2}	10^{-1}
$GTA5 \rightarrow CS$		33.6	34.6	38.2	34.6
$CS \rightarrow Mapillary$	42.6	43.0	45.0	37.6	

Table 4. $GTA5 \rightarrow Cityscapes$. First row, source, indicates a model trained on GTA5 only. We also report performance on two unseen domains, Mapillary and Apolloscape. All units are % mIOU.

		$GTA5 \rightarrow Cityscapes$			
Method		Tested on			
		GTA5 (Source)	Cityscapes (Target)	Mapillary (Unseen)	Apolloscape (Unseen)
Source		62.2	34.8	37.1	25.3
UADA		62.7	38.2	38.5	27.4

5.2. Unsupervised adversarial domain adaptation

In this experiment, we evaluate UADA using DAN. We consider two adaptation tasks, compare with related works and report results on unseen domains, Mapillary and Apolloscape. We found the values for λ using a hyperparameter sweep and report results in Table 3.

Table 4 reports the results for adapting from GTA5 to

Table 5. **Comparison with three recent works in UADA.** The first column indicates at what level the adversary operates. Image size indicates number of pixels in height \times width. Performance units are % mIOU.

$GTA5 \rightarrow Cityscapes$			
Adversarial adaptation level	Norm. layer	Image size	Target perf.
Representation [19]	No	x	25.5
Representation and logit [47]	No	x	28.9
Representation and output [35]	No	512×1024	37.1
Representation [ours]	Yes	384×768	38.2

Cityscapes. These results show that we improve 3.4 points mIOU on the target domain. A source-only model achieves 62.1 % mIOU on GTA5 and 34.8 % mIOU on Cityscapes. UADA achieves 62.7 % mIOU on GTA5 and 38.2 % mIOU on Cityscapes. Referring back to Figure 4, we observe that the retrieval curve for the adapted model (yellow) lies closer to the perfect alignment curve, which confirms that UADA aligns the representations. We plot segmentation outputs of both models on the Cityscapes validation set in Figure 5. We observe that the source-only model, trained on GTA5, makes incorrect predictions in large areas of the Cityscapes images. The plots for UADA clearly show less of these er-

Table 6. *Cityscapes* \rightarrow *Mapillary*. The first row, source, indicates a model trained on Cityscapes only. We also report performance on an unseen domain, Apolloscape. All units are % mIOU.

<i>Cityscapes</i> \rightarrow <i>Mapillary</i>			
Method	Cityscapes (Source)	Tested on Mapillary (Target)	Apolloscape (Unseen)
Source	63.6	43.2	25.8
UADA	64.0	45.0	27.1

rors. In Table 5, we compare our results to recent publications on this adaptation task. Despite using lower resolution, training UADA using DAN achieves the highest performance, 38.2 % mIOU, on the target domain.

Table 4 shows that the performance improves on two unseen domains, Mapillary and Apolloscape. For Mapillary, the source-only model achieves 37.1 % mIOU and UADA improves the performance to 38.5 % mIOU. For Apolloscape, the source-only model achieves 25.3 % mIOU and UADA improves the performance to 27.4 % mIOU. These results show that by adapting from GTA5 to Cityscapes, the performance improves 1.4 points mIOU on Mapillary and 2.1 points mIOU on Apollo. Appendix 5 shows segmentation outputs of both models on an unseen domain, Mapillary. Again on this unseen domain, we observe less incorrect predictions when comparing the unadapted with the adapted model.

Finally, we reproduce the benefits of DAN on a second, real world adaptation task: from Cityscapes to Mapillary. Both domains are real and have less domain gap. Consequently, we expect a smaller improvement compared to synthetic to real adaptation. Table 6 reports the results for this adaptation task. We observe that the source domain improves from 63.6 to 64.0 % mIOU. The target domain improves 1.8 points from 43.2 to 45.0 % mIOU. For this second adaptation task, we evaluate the source-only model and the adapted model on the unseen domain, Apolloscape. The performance on Apolloscape improves 1.3 points from 25.8 to 27.1 % mIOU. Although these improvements are less than adapting from a synthetic to a real domain, this experiment again shows the improved generalization capability for unseen domains, which is a chief aim in pattern recognition.

6. Discussion

Our experimental results demonstrate that UADA with our normalization layer improves segmentation performance on source, target and unseen domains. In this section, we highlight some points for discussion.

In training UADA, the λ hyperparameter plays a critical role for the performance on the target domain. The λ bal-

ances the segmentation and confusion loss, see Equation 2. In Table 3, we report the target performance on both adaptation tasks using different values for λ . We observe that 1) the performance drops more than 2 points mIOU when we change λ from the optimal value; 2) the optimal value differs between the adaptation tasks. However, in unsupervised domain adaptation, we would have no labels for the target domain to evaluate these numbers. One might thus argue that tuning the hyperparameters alludes to overfitting on the target labels. Therefore, we evaluate the model on unseen domains, i.e. the domains that we never used training nor setting the hyperparameters. For both adaptations tasks in Tables 4 and 6, we observe an improvement in the performance on unseen domains. These improvements show that training with the optimal λ value for the target domain also improves performance for unseen domains.

The domain classifier operates at the level of the representations, which might not suffice to reduce the domain adaptation gap. An adapted model from GTA5 to Cityscapes achieves 38.2 % mIOU, while a model trained and evaluated on Cityscapes achieves 63.6 % mIOU. This gap shows the need for further improvement. The first column in Table 5 shows at what level the adversaries operate. Contemporary work focuses on extending the adversaries to multiple levels, [35, 10]. These works show a promising direction of research and we believe that our DAN will facilitate further improvements.

7. Conclusion

In this work, we present a Domain Agnostic Normalization (DAN) layer for unsupervised adversarial domain adaptation (UADA). Unlike conventional batch normalization, for which we show that it significantly reduces performance on UADA, our DAN layer learns normalization statistics only over the labeled domain and applies it to unlabeled and unseen domains. Consequently, DAN unlocks the benefits of normalization for UADA and we present in-depth experiments to demonstrate its advantage. Despite training at less than half of the original resolution, we surpass state-of-the-art performance on a common benchmark, i.e. adapting from GTA5 to Cityscapes, achieving 38.2 % mIOU on the target domain. As a model trained with DAN remains domain agnostic when testing, we also report 1.4 and 2.1 points mIOU improvement on two unseen domains, Mapillary and Apolloscape, respectively. We reproduce these improvements on a second adaptation task, i.e. adapting from Cityscapes to Mapillary. Again, we show a performance improvement on an unseen domain, Apolloscape. All together, we demonstrate that using DAN allows for improved performance on unseen domains, which is the chief aim of pattern recognition.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, L. Kaiser, M. Kudlur, J. Levenberg, D. Man, R. Monga, S. Moore, D. Murray, J. Shlens, B. Steiner, I. Sutskever, P. Tucker, V. Vanhoucke, V. Vasudevan, O. Vinyals, P. Warden, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *OSDI*, 16:265—283, 2016.
- [2] M. Arjovsky and L. Bottou. Towards principled methods for training generative adversarial networks. *International Conference on Learning Representations (ICLR)*, 2017.
- [3] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv*, 2016.
- [4] V. Badrinarayanan, A. Kendall, and R. Cipolla. SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481—2495, 2017.
- [5] S.-H. Bae, Y. Lee, Y. Jo, Y. Bae, and J.-w. Hwang. Rank of Experts: Detection Network Ensemble. *arXiv*, 2017.
- [6] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, J. Wortman Vaughan, S. R. Ben-David David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. Vaughan. A theory of learning from different domains. *Machine Learning*, 79(79):151–175, 2010.
- [7] Y. Bengio, A. Courville, and P. Vincent. Representation Learning: A Review and New Perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, aug 2013.
- [8] J. Bjorck, C. Gomes, and B. Selman. Understanding Batch Normalization. *arXiv*, 2018.
- [9] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *IEEE Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 7, 2017.
- [10] Y.-H. Chen, W.-Y. Chen, Y.-T. Chen, B.-C. Tsai, Y.-C. Wang, and M. Sun. No More Discrimination: Cross City Adaptation of Road Scene Segmenters. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2011—2020, apr 2017.
- [11] S. Chintala. How to train a GAN, Workshop on Adversarial Training at NIPS 2016, 2016.
- [12] 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 *IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 3213—3223, apr 2016.
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 248–255. IEEE, 2009.
- [14] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1915–1929, 2013.
- [15] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-Adversarial Training of Neural Networks. *Journal of Machine Learning Research (JMLR)*, 17(1):189–209, 2017.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative Adversarial Nets. *Advances in Neural Information Processing Systems (NIPS)*, pages 2672–2680, 2014.
- [17] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In *European Conference on Computer Vision (ECCV)*, pages 297–312. Springer, 2014.
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [19] J. Hoffman, D. Wang, F. Yu, and T. Darrell. FCNs in the Wild: Pixel-level Adversarial and Constraint-based Adaptation. *ArXiv*, dec 2016.
- [20] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. *arXiv*, 2017.
- [21] X. Huang, X. Cheng, Q. Geng, B. Cao, D. Zhou, P. Wang, Y. Lin, and R. Yang. The ApolloScape Dataset for Autonomous Driving. *arXiv*, 2018.
- [22] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, pages 448–456, 2015.
- [23] A. Khosla, T. Zhou, T. Malisiewicz, A. A. Efros, and A. Torralba. Undoing the damage of dataset bias. In *European Conference on Computer vision (ECCV)*, volume 7572, pages 158–171, 2012.
- [24] J. Kohler, H. Daneshmand, A. Lucchi, M. Zhou, K. Neymeyr, and T. Hofmann. Towards a Theoretical Understanding of Batch Normalization. *ArXiv*, 2018.
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1097–1105, 2012.
- [26] Y. Li, N. Wang, J. Shi, J. Liu, and X. Hou. Revisiting Batch Normalization For Practical Domain Adaptation. In *International Conference on Learning Representations (ICLR)*, page Workshop track, 2017.
- [27] Z. Liang, Y. Feng, Y. Guo, H. Liu, L. Qiao, W. Chen, L. Zhou, and J. Zhang. Learning for Disparity Estimation through Feature Constancy. *arXiv*, 2017.
- [28] M.-Y. Liu and O. Tuzel. Coupled Generative Adversarial Networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems (NIPS)*, pages 469–477. Curran Associates, Inc., 2016.
- [29] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 3431–3440, 2015.
- [30] G. Neuhold, T. Ollmann, S. R. Bulò, and P. Kotschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *International Conference on Computer Vision (ICCV)*, pages 22–29, 2017.

- [31] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1520–1528, 2015.
- [32] D. Pathak, E. Shelhamer, J. Long, and T. Darrell. Fully Convolutional Multi-Class Multiple Instance Learning. *arXiv*, 2014.
- [33] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. In European Conference on Computer vision (ECCV), editor, *Lecture Notes in Computer Science*, volume 9906, pages 102–118. Springer, Cham, oct 2016.
- [34] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2234–2242, 2016.
- [35] S. Sankaranarayanan, Y. Balaji, A. Jain, S. Lim, and R. Chellappa. Unsupervised Domain Adaptation for Semantic Segmentation with GANs. *ArXiv*, nov 2017.
- [36] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry. How Does Batch Normalization Help Optimization? (No, It Is Not About Internal Covariate Shift). *ArXiv*, 2018.
- [37] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv*, 2014.
- [38] T. Tommasi, N. Patricia, B. Caputo, and T. Tuytelaars. A deeper look at dataset bias. In *German conference on Pattern Recognition (GCPR)*, pages 504–516, 2015.
- [39] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 1521–1528. IEEE, jun 2011.
- [40] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous Deep Transfer Across Domains and Tasks. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4068—4076, dec 2015.
- [41] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial Discriminative Domain Adaptation. *ArXiv*, 2017.
- [42] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance Normalization: The Missing Ingredient for Fast Stylization. *ArXiv*, 2016.
- [43] U. Unknown. Removed for blind review. 2018.
- [44] S. Xiang and H. Li. On the Effects of Batch and Weight Normalization in Generative Adversarial Networks. *stat*, 1050:22, 2017.
- [45] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 5987–5995. IEEE, 2017.
- [46] F. Yu and V. Koltun. Multi-Scale Context Aggregation by Dilated Convolutions. *arXiv*, 2015.
- [47] Y. Zhang, P. David, and B. Gong. Curriculum Domain Adaptation for Semantic Segmentation of Urban Scenes. In *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [48] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired Image-To-Image Translation Using Cycle-Consistent Adversarial Networks. In *IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 2223–2232, 2017.