

Learning to adapt class-specific features across domains for semantic segmentation

Mikel Menta*, Adriana Romero†, Joost van de Weijer*

*Computer Vision Center, Universitat Autònoma de Barcelona

†McGill University, Montreal

Abstract—Recent advances in unsupervised domain adaptation have shown the effectiveness of adversarial training to adapt features across domains, endowing neural networks with the capability of being tested on a target domain without requiring any training annotations in this domain. The great majority of existing domain adaptation models rely on image translation networks, which often contain a huge amount of domain-specific parameters. Additionally, the feature adaptation step often happens globally, at a coarse level, hindering its applicability to tasks such as semantic segmentation, where details are of crucial importance to provide sharp results. In this thesis, we present a novel architecture, which learns to adapt features across domains by taking into account per class information. To that aim, we design a conditional pixel-wise discriminator network, whose output is conditioned on the segmentation masks. Moreover, following recent advances in image translation, we adopt the recently introduced *StarGAN* architecture as image translation backbone, since it is able to perform translations across multiple domains by means of a single generator network. Preliminary results on a segmentation task designed to assess the effectiveness of the proposed approach highlight the potential of the model, improving upon strong baselines and alternative designs.

Index Terms—domain adaptation, semantic segmentation, adversarial training, conditional discriminator, convolutional neural networks

I. INTRODUCTION

IN the last decade, deep learning has become the *de facto* standard in many machine learning application domains such as computer vision, natural language processing or speech. More specifically, in computer vision, the success of AlexNet [29] in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2012 changed the research landscape. Since then, deep learning architectures have been quickly spreading and have shown impressive results in tasks such as image classification [21], [23], [29], [48], [50], semantic segmentation [3], [7], [8], [13], [26], [43], [46], [56], [59], or object detection [16], [40], [41], among many others.

These recent successes of deep neural networks in many application domains have been attributed, within significant degree, to the availability of large-scale labeled datasets that can be used for training, the increase of computational power through GPUs and the development of sophisticated approaches. However, even if these networks are generally able to achieve a high performance in unseen data samples from the same dataset, they often don't generalize well to new datasets and tasks. This problem is referred to as domain shift or bias. A typical procedure to tackle the domain shift across different datasets is to fine-tune the network trained on a dataset, with samples from the new dataset or task. However, this procedure usually requires having access to a considerable amount of additional annotated data.

Furthermore, tasks such as semantic segmentation especially suffer from an extremely expensive annotation, due to the pixel-wise nature of the task. To mitigate labeling efforts, many works suggest exploiting datasets built from data generated by computer simulated environment, such as GTA5 [42] and SYNTHIA [44], where infinite

amounts of data can be easily collected. Unfortunately, training on simulated data leads to the above-mentioned domain shift, when applying the trained models to a real environment, achieving poor generalization.

Motivated to overcome the domain shift problem, domain adaptation methods propose a handful of approaches to improve the knowledge transfer from a source domain to a target domain. Nevertheless, in the extreme setting of unsupervised domain adaptation, the knowledge needs to be adapted without any target domain annotation. Recent advances in unsupervised domain adaptation mostly target the task of image classification [14], [15], [34], [47], [52], with several extensions to other tasks such as semantic segmentation [22], [37], [51] and object detection [9], [24].

The majority of the recent domain adaptation literature tries to obtain invariant representations between both source and target domains through adversarial training, following [14]. In this adversarial game, generally, the discriminator network aims to distinguish the original domain of the feature representation, while the feature encoding network intends to extract indistinguishable representations from both domains by fooling the discriminator. It is worth mentioning that this adversarial cross-domain feature matching is broadly performed at a coarse level, in a global way [9], [22], [24], [37], [51], without paying special attention to fine-grained information such as object positions/sizes or, more generally, class distributions, which are relevant for tasks such as semantic segmentation. Note that extracting such detailed information increases the challenges of domain adaptation per se.

Additionally, many approaches incorporate image translation or reconstruction objective to obtain more general features, which may be composed of separate networks to process each domain. These auxiliary objectives often introduce additional modules, leading to a vast amount of parameters, which are then discarded during the testing phase.

In this thesis, we present a novel unsupervised domain adaptation model for semantic segmentation, which addresses the previous concerns. The contributions are twofold:

- 1) We successfully reduce the amount of parameters of our domain adaptation network to roughly the half w.r.t. many other methods, by following the idea of *StarGAN* of providing an additional channel to the input of the network with the corresponding domain label.
- 2) We present a novel pixel-wise discriminator training procedure that is able to perform a local discrimination of the feature representations by taking into account per class information. This is achieved by applying the segmentation predictions to the discriminator output, allowing the network to perform the feature distribution matching at class level, instead of globally.

To assess the impact of the contributions of this thesis, we modify MNIST dataset to simulate characteristics that we encounter in real scenarios, namely differences in appearance (texture, color, brightness) between different domains, as well as in per class pixel distributions. We report both qualitative and quantitative results in

those MNIST modifications as a proof of concept of the method proposed in this thesis and show that our local class specific feature matching significantly outperforms the commonly employed global feature matching, as well as the baselines. We also provide an ablation study of our model from the number of parameters in the model's backbone perspective, in addition to the adversarial feature matching, outlining the benefits of our contribution.

The rest of the thesis is organized as follows. Section II reviews the recent state-of-the-art literature. Then, Section III introduces the concepts needed for the understanding of this thesis. After that, Section IV details the proposed approach. Section V motivates the experiments and comparisons w.r.t. baselines and ablated models, Section VI reports the obtained qualitative and quantitative results, and finally, Section VII draws the conclusions of this thesis and presents potential future research directions.

II. STATE OF THE ART

This section aims to give the reader an overview of recent state-of-the-art approaches that tackle problems relevant to this thesis, namely image translation, semantic segmentation and unsupervised domain adaptation, with special focus on the latter.

A. Image translation

Image translation and style transfer consist on transferring the appearance of a certain reference image or domain of images to a target image. These areas have been benefiting from a large amount of contributions in recent years. In works like [27], [53]–[55], approaches are based on information extracted from increasingly abstract layers of a pretrained network. The goal can be summarized as trying to obtain a low level style representation similar to the reference image, while preserving more abstract content-based features.

Nevertheless, with the introduction of Generative Adversarial Networks (GANs) [18] and conditional GANs [36], new research directions exploiting and extending these architectures to perform image translation have emerged [5], [25], [60]. Although adversarial architectures have been enjoying increasing popularity among researchers, alternatives based on the combination of features extracted from encoder-decoder architectures have also shown competitive results [17].

A possible way of categorizing image translation approaches is by dividing them into: paired case, where for each input image a corresponding one in the opposite domain is provided; and unpaired case, where there is not such correspondence between images of both domains.

In the paired case, we find works such as [17], [25], among others. On one hand, [25] proposed to generate the translated image with a GAN conditioned on an input image, guiding the network with its corresponding paired image. On the other hand, [17] tries to obtain a disentangled representation of the image with both a domain agnostic and a domain specific part. To do so, they train an encoder and a decoder per domain using an aggregation of several loss components.

In the unpaired case, approaches are slightly different. For example, the solution provided in [60] (explained with more details in Section III-B) is to learn translation mappings among domains (from domain A to domain B and viceversa) by means of adversarial training, and by introducing a cycle-consistency to enforce that, e.g. the translation of an image from A to B, followed by the translation from B to A, leads back to the original image. Similarly, in [10] they propose to use a single generator network, together with a discriminator, to perform translations across multiple domains at the expense of adding a few conditional parameters (detailed in Section III-C). Finally, in [30] the strategy is focused on obtaining

representative features that are agnostic to the domains and that, as a result, can be decoded to any target domain

B. Semantic segmentation

The task of semantic segmentation, which consists on classifying each pixel of an image according to some known labels, is currently mainly addressed by Deep Neural Networks (DNN) in the state-of-the-art. Current DNN models to tackle pixel-prediction problems are based on Fully Convolutional Networks (FCNs) [46]. FCNs endow Convolutional Neural Networks (CNNs) with an upsampling path to recover the input resolution. These networks can be trained end-to-end.

Some of the contributions in this area have focused on obtaining structurally better representations by keeping the pooling indices of the downsampling steps to perform a better upsampling [3], adding skip connections to preserve low level spatial features [13], [43] or combining features at different spatial resolutions [7], [26], [59]. Other works have tried to overcome the loss of spatial resolution introduced by pooling layers by means of dilated convolutions to enlarge the receptive field [8], [56] while others have analyzed the benefits and limitations of different receptive field enlargement operations [7].

Unluckily, an important issue of the semantic segmentation task is the high cost of obtaining labeled data and the big amounts of data required to train high capacity deep learning models. A proposed solution for this problem has been to use weak annotations, which can be obtained in an easier (and less expensive) manner [4], [12], [32]. Other proposed alternatives involve training the models on synthetic datasets such as SYNTHIA [44] or GTA5 [42], where labels can be extracted automatically. Unfortunately, models trained on synthetic datasets tend to not generalize well to real world scenarios.

C. Unsupervised domain adaptation

Domain adaptation methods aim to address the challenges posed by existing dataset shifts (e.g. when training and testing samples come from different distributions) with the goal of transferring the knowledge learned in a source domain to a target domain. A well known approach for this problem is training in the source domain and *fine-tuning* in the target domain. However, in many real cases, the ground truth information for the target domain might not be available (*unsupervised domain adaptation*) or might be too sparse (*semi-supervised domain adaptation*).

Due to the usefulness of solving the unsupervised case, we can find a handful of computer vision work in this area. The majority of them focus or experiment with the *classification* task [6], [14], [15], [20], [33], [34], [34], [38], [45], [47], [52], [57], [61]. But we can also find a significant variety of approaches in semantic segmentation [22], [37], [51], object detection [9], [24] and depth estimation [2], [5], among others.

The methodologies developed to tackle domain adaptation can vary on some specific aspects depending on the particularities of the task. However, a significant number of works rely on state-of-the-art image translation methods (reviewed in subsection II-A) or simple image reconstruction with different purposes. Some of these purposes are: (1) obtaining unsupervised features, reconstructing the input images, to perform the task directly [15]; (2) translating source images to the target domain and training a network from those images in a supervised manner [2], [5] or (3) translating the source images and using them in a subsequent fine-tuning step [24]. Generally, in many cases, the image translation and/or reconstruction tasks are used as additional signal to learn more representative features for the target domain [6], [37].

It is worth mentioning that the vast majority of approaches to (unsupervised) domain adaptation focus on *matching distributions* of source and target domains at feature level. To do so, state-of-the-art approaches rely on (1) maximizing the correlation between features of both domains [61], (2) associating the target samples with the source samples [20] or (3) obtaining a clustered distribution in the features of the target domain [47]. Nonetheless, a vast set of the recent literature focuses on matching distributions of source and target domains via *adversarial matching* of the features, either by means of a domain discriminating classifier, a *Gradient Reversal Layer* (first introduced in [14]) or the vanilla adversarial objective introduced in [18]. Some of these approaches propose to learn translated features from source to target domains following the image translation paradigm [22]. Moreover, the feature matching can also be done at different depth levels of the network simultaneously [34], [51]. A complementary view is to isolate the domain agnostic features from the domain specific ones as in [6] using different network modules. In [52] the authors analyze the previously mentioned alternatives along with an overview of general challenges encountered in domain adaptation. Similarly, in [37] the authors exploit and combine several of the previously mentioned components into a single model to obtain better training signals. Recently, works on modeling the domain adaptation problem as a similarity learning task have also emerged in the literature [38]. Finally, alternative research directions include (1) predicting the weights of the target model [45]; (2) addressing the well known mode collapse problem in adversarial training with conditioning [33]; (3) dealing with a partial unsupervised domain adaptation formulation [57]; or (4) addressing and analyzing less popular tasks in the domain adaptation literature, such as object detection [9].

III. BACKGROUND

In this section, we will go through the design and characteristics of some models needed for the understanding of this thesis. In the following subsections, we will explain fully convolutional network architectures for semantic segmentation [13], [26], [43], [46] as well as the *CycleGAN* [60] and *StarGAN* [10] architectures for unpaired image translation.

A. Fully convolutional network for semantic segmentation

As mentioned in Section II-B, a common approach for tackling the semantic segmentation problem is the use of FCNs that allow to obtain a segmentation mask directly from an image in an end-to-end way. FCNs are composed of a downsampling path (with convolutional and subsampling operations) followed by an upsampling path (with transposed convolutions), which recovers the input resolution. In this thesis, we also use an FCN for performing the segmentation task.

Given a dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$ where $(x^{(i)}, y^{(i)})$ are pairs of image and segmentation masks respectively, we define a FCN S which approximates as well as possible the mapping of an input sample x , potentially not included in \mathcal{D} , to its corresponding label y . Figure 1 depicts an FCN architecture with an initial convolutional block, 2 downsampling blocks, followed by 2 upsampling blocks and a final convolution that performs the segmentation. Extra details of the implementation and training are explained in Appendix A. As will be later explained in the thesis, this model is used as a baseline (Section V-B) and its decoder, composed of the upsampling path exclusively, is used as segmenter module in our model (Section IV).

In the binary segmentation case, we can train S by minimizing the *soft IOU loss* [13], [39] between the predictions $\hat{y} = S(x)$ and ground-truth segmentations y . Considering binary segmentations, where y^p is 1 if pixel p in segmentation y corresponds to foreground

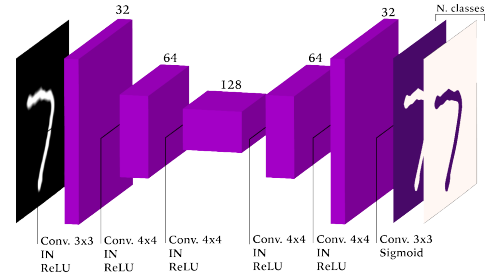


Figure 1. Architecture of the fully convolutional network for semantic segmentation defined in Section III-A. *IN* indicates instance normalization.

class and 0 otherwise, the segmentation loss \mathcal{L}_{segm} is defined as follows:

$$\mathcal{L}_{segm} = \mathbb{E}_x \left[1 - \frac{\sum_p \hat{y}^p \cdot y^p}{\sum_p (\hat{y}^p + y^p - \hat{y}^p \cdot y^p)} \right]. \quad (1)$$

Note that for multi-class segmentation problems, the cross-entropy loss is used as de facto standard. Considering the segmentations in a *one-hot* encoding, where $y^{c,p}$ is 1 if pixel p in segmentation y corresponds to class c and 0 otherwise, the cross-entropy loss is defined as follows:

$$\mathcal{L}_{segm} = \mathbb{E}_{x,p,c} [-y^{c,p} \cdot \log(\hat{y}^{c,p}) - (1 - y^{c,p}) \cdot \log(1 - \hat{y}^{c,p})]. \quad (2)$$

B. CycleGAN

One of the most popular image translation architectures based in conditional adversarial models is *CycleGAN* [60]. *CycleGAN* is designed to perform unpaired image translations between two domains of images by training two FCN generators and two discriminators. Each generator is trained to translate an image from one domain to the other.

Formally, given two datasets of images $\mathcal{D}_a = \{x_a^{(i)}\}_{i=1}^{N_a}$ and $\mathcal{D}_b = \{x_b^{(i)}\}_{i=1}^{N_b}$ sampled from two different domains a and b , we define two generator networks $G_{ab} : x_a \rightarrow x_b$ and $G_{ba} : x_b \rightarrow x_a$. Apart from those, we define a discriminator per domain \mathcal{D}_a and \mathcal{D}_b which aim to distinguish images sampled from their corresponding dataset from the ones generated by G_{ab} and G_{ba} . The generator networks are FCNs which, in our implementation, have an initial convolutional block, followed by 2 downsampling blocks, 2 upsampling blocks and a final convolution that maps back to the image domain. The discriminators D follow the *PatchGAN* architecture [25] with 4 downsampling blocks and a fully connected layer with a single output. These networks are depicted in Figure 2 and overall, they are trained following two objectives:

Adversarial losses. In order to obtain realistic translations for each domain, an adversarial loss is optimized between each pair of generator and discriminator networks. Following the traditional adversarial setting, the discriminator tries to distinguish if an image comes from a real dataset or has been generated, whereas the generator tries to fool the discriminator. Under this idea and, using the Least Squares GAN loss [35], the following loss components are defined for images from domain a :

$$\mathcal{L}_{adv}^{D_a} = \mathbb{E}_{x_a} [(D_a(x_a) - 1)^2] + \mathbb{E}_{x_b} [D_a(G_{ba}(x_b))^2] \quad (3)$$

$$\mathcal{L}_{adv}^{G_{ba}} = \mathbb{E}_{x_b} [(D_a(G_{ba}(x_b)) - 1)^2] \quad (4)$$

and analogously, the following are defined for images from domain b :

$$\mathcal{L}_{adv}^{D_b} = \mathbb{E}_{x_b} [(D_b(x_b) - 1)^2] + \mathbb{E}_{x_a} [D_b(G_{ab}(x_a))^2] \quad (5)$$

$$\mathcal{L}_{adv}^{G_{ab}} = \mathbb{E}_{x_a} [(D_b(G_{ab}(x_a)) - 1)^2] \quad (6)$$

The setting of the subnetworks to compute the above-mentioned loss components are represented in Figure 2.

Cycle consistency loss. The adversarial loss components don't encourage the model to keep the content intact in the generated translations, they only focus on producing images that look realistic in the translated domain. In order to motivate this behaviour, a cycle consistency loss is added, so that an image translated to the opposite domain and back to the original domain remains invariant. The cycle consistency loss is defined as:

$$\mathcal{L}_{cyc} = \mathbb{E}_{x_a}[\|G_{ba}(G_{ab}(x_a)) - x_a\|_1] + \mathbb{E}_{x_b}[\|G_{ab}(G_{ba}(x_b)) - x_b\|_1]. \quad (7)$$

An example of the network setting in one of the two possible cycles is represented in Figure 3. Note that an analogous pipeline is used to compute the cycle loss for the remaining domain.

Finally, the model is trained by making the generators G_{ab} and G_{ba} minimize the objective

$$\mathcal{L}^G = \mathcal{L}_{adv}^{G_{ab}} + \mathcal{L}_{adv}^{G_{ba}} + \lambda \cdot \mathcal{L}_{cyc}, \quad (8)$$

where λ is an hyper-parameter controlling the relative importance of the two objectives, and making the discriminators D_a and D_b minimize the objective

$$\mathcal{L}^D = \mathcal{L}_{adv}^{D_a} + \mathcal{L}_{adv}^{D_b}. \quad (9)$$

C. StarGAN

StarGAN [10] is a recently proposed alternative to *CycleGAN* to address image translation, and is also the model that we choose as backbone for the unsupervised domain adaptation pipeline introduced in Section IV. *StarGAN* offers several practical advantages w.r.t. *CycleGAN*, namely it is designed to perform translations across multiple domains and uses a single generator network and a single discriminator network. To do so, the model expects to receive the translation domain label together with the image to be translated as input to the generator network, performing a different translation according to this label.

Given that in the problem of domain adaptation only a source and a target domain are taken into account, we simplify the notation of the equations to the case of having only two domains. So, considering the datasets \mathcal{D}_a and \mathcal{D}_b introduced in Section III-B, we define a generator network G that, given an image $x \in \{\mathcal{D}_a \cup \mathcal{D}_b\}$ with its corresponding domain label $l \in \{a, b\}$ and opposite label $\bar{l} \in \{a, b\}$ with $\bar{l} \neq l$, is able to perform a translation $\hat{x}_{\bar{l}} = G(x_l, \bar{l})$. The translation label \bar{l} is replicated along the spatial dimension to match the size of x . We also define a discriminator D that, given an image x , predicts (a) whether x is sampled from one of the datasets or has been generated by G (we will denote this prediction as $D_{rf}(x)$) and (b) the domain label of x (we will denote this prediction as $D_{dom}(x)$).

Similarly to the *CycleGAN* model explained in Section III-B, the generator G is a FCN, which in our implementation has an initial convolutional block, followed by 2 downsampling blocks, 2 upsampling blocks and a final convolution that maps back to the image domain. Again, as in *CycleGAN*, the discriminator D follows the *PatchGAN* [25] discriminator architecture with 4 downsampling blocks and two outputs D_{rf} and D_{dom} . The network is shown in Figure 4.

In order to train these modules, three different objectives are defined following similar goals to the ones *CycleGAN*:

Translation adversarial loss. In order to make the generated image-translations look realistic, an adversarial objective is used

between G and D_{rf} , implemented as the Wasserstein loss with gradient penalty [19].

$$\mathcal{L}_{rf}^D = \mathbb{E}_{x, \bar{l}}[D_{rf}(G(x, \bar{l}))] - \mathbb{E}_x[D_{rf}(x)] + \lambda_{gp} \mathbb{E}_{\hat{x}}[(\|\nabla_{\hat{x}} D_{rf}(\hat{x})\|_2 - 1)^2], \quad (10)$$

$$\mathcal{L}_{rf}^G = -\mathbb{E}_{x, \bar{l}}[D_{rf}(G(x, \bar{l}))], \quad (11)$$

where \hat{x} is sampled uniformly along a straight line between a pair of real and generated images and λ_{gp} controls the relative importance of the *gradient penalty* component. The setting of generator and discriminator networks to compute these loss components is represented in Figure 4.

Domain classification loss. The goal of translation is to have the transformed images look from the translated domain, which can also be interpreted as having them classified as belonging to the translation domain. To do so, D_{dom} is trained to classify real images as belonging to their corresponding domain and G is trained to make the translated images correctly classified by D_{dom} . This training is performed by minimizing

$$\mathcal{L}_{dom}^D = \mathbb{E}_{x_a}[H(D_{dom}(x_a), a)] + \mathbb{E}_{x_b}[H(D_{dom}(x_b), b)], \quad (12)$$

$$\mathcal{L}_{dom}^G = \mathbb{E}_{x_a}[H(D_{dom}(G(x_a, b)), b)] + \mathbb{E}_{x_b}[H(D_{dom}(G(x_b, a)), a)], \quad (13)$$

where H denotes the *cross-entropy* loss defined as $H(x, y) = -y \cdot \log(x) - (1 - y) \cdot \log(1 - x)$. The setting of the generator and discriminator networks to compute these loss components is analogous to the translation adversarial component shown in Figure 4.

Cycle consistency loss. In order to preserve the content of the input image in the translations, a cycle consistency loss is added to the generator's set of losses, following *CycleGAN*:

$$\mathcal{L}_{cyc} = \mathbb{E}_{x_a}[\|G(G(x_a, b), a) - x_a\|_1] + \mathbb{E}_{x_b}[\|G(G(x_b, a), b) - x_b\|_1], \quad (14)$$

An example of one of the two cycle settings is represented in Figure 5.

To wrap up, the model is trained by minimizing \mathcal{L}^G for the generator G and \mathcal{L}^D for the discriminator D :

$$\mathcal{L}^G = \lambda_{rf} \cdot \mathcal{L}_{rf}^G + \lambda_{dom} \cdot \mathcal{L}_{dom}^G + \lambda_{cyc} \cdot \mathcal{L}_{cyc} \quad (15)$$

$$\mathcal{L}^D = \lambda_{rf} \cdot \mathcal{L}_{rf}^D + \lambda_{dom} \cdot \mathcal{L}_{dom}^D \quad (16)$$

where λ_{rf} , λ_{dom} and λ_{cyc} control the relative importance of the three objectives.

IV. METHOD

In the setting of unsupervised domain adaptation, we assume access to a dataset $\mathcal{D}_s = \{(x_s^{(i)}, y_s^{(i)})\}_{i=1}^{N_s}$ drawn from a source domain distribution, where in the case of image segmentation, x denote images and y ground-truth segmentation masks. At the same time, we assume access to another dataset $\mathcal{D}_t = \{(x_t^{(n)})\}_{n=1}^{N_t}$ sampled from a target domain distribution. Note that in \mathcal{D}_t the ground-truth segmentation masks are not provided. With these two datasets the goal is to learn a function $f : x \rightarrow y$ that correctly predicts both y_s and y_t given x_s and x_t , respectively and generalizes properly to unseen source and target samples.

In order to tackle the unsupervised domain adaptation problem, we introduce a model based on an image translation backbone that we use to translate images x_s and x_t to their opposite domain. Then, a segmenter network is connected to the bottleneck features of the image translator to perform our end task (semantic segmentation). The bottleneck features of the image translator are also connected

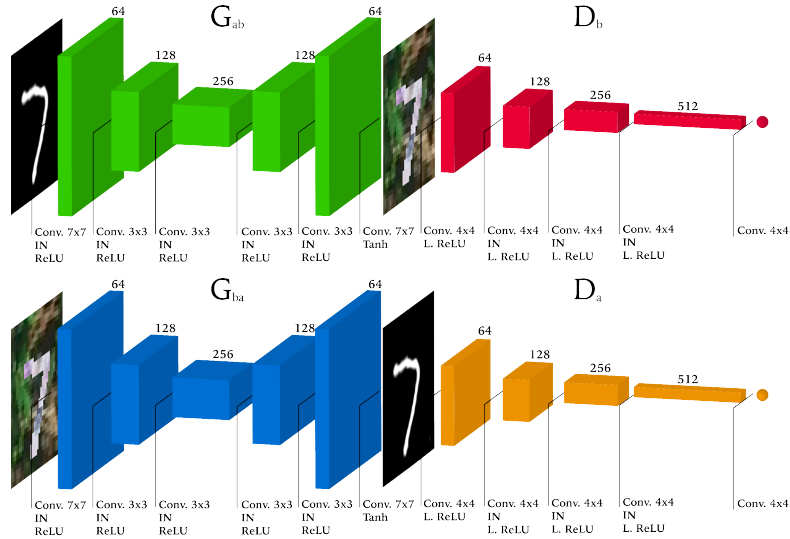


Figure 2. Architecture of the *CycleGAN* [60] presented in Section III-B. *IN* indicates instance normalization and *L. ReLU* leaky rectified linear unit.

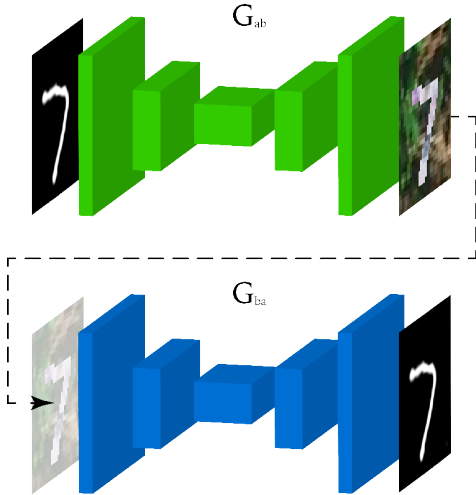


Figure 3. Cycle configuration of the modules in *CycleGAN* [60] in order to compute \mathcal{L}_{cyc} for domain a .

to an additional module, which performs a class specific feature matching. The rationale behind this feature matching module is to ensure that features belonging to the same class come from the same distribution. An overview of the model is depicted on Figure 6: note that G_e and G_d depict the image translation network, followed by a discriminator D ; S refers to the segmenter; and D_f denotes the class-aware feature matching module. The rest of the section will be devoted to providing details on each component of the proposed model.

Image translation backbone. As mentioned in Section II-C, many domain adaptation architectures use image translation or image reconstruction as an auxiliary task to obtain more general features, in order to boost the classification or segmentation performance on the target domain. In our model, we use image translation as auxiliary task and we rely on the state-of-the-art *StarGAN* [10] model as the backbone module of the proposed network.

As already explained in Section III-C, *StarGAN* is able to perform unpaired image translations between a set of different image domains, so we are able to exploit it to translate our images from \mathcal{D}_s (source domain) to \mathcal{D}_t (target domain) and the other way around. The main

advantage of using *StarGAN* architecture is that it allows us to have domain specific parameters in a single network, drastically reducing the number of parameters required to perform image translation w.r.t. related approaches in the literature such as [2], [6], [34], [47], [52].

Therefore, our model's base is composed of both *StarGAN* generator G and discriminator D networks. These networks are trained with the *StarGAN* objectives introduced in Section III-C. Using x_s and x_t as input images, we define:

- \mathcal{L}_{rf}^G and \mathcal{L}_{rf}^D (Eqs. (11) and (10)) to encourage the model to generate realistic looking translations.
- \mathcal{L}_{dom}^G and \mathcal{L}_{dom}^D (Eqs. (13) and (12)) to make the translated images look like drawn from their corresponding translation domain.
- \mathcal{L}_{cyc} (Eq. (14)) to encourage the network to preserve the content of the images while translating them.

Segmenter. To perform the segmentation task, we connect a segmenter decoder to the bottleneck of *StarGAN*. To do so, we divide the generator G into the encoder composed of the downsampling path of G and, denoted from now on as G_e , and the decoder composed of the upsampling path of G and referred to as G_d . This way, we connect a segmentation decoder S to G_e , and feed both G_d and S with the output of G_e . Note that the segmentation decoder S is built following the upsampling path of the segmenter defined in Section III-A. The above-described connection between G and S is depicted in Figure 6.

The segmentation decoder S is trained by minimizing the *soft IOU loss* defined in Eq. (1) between y_s and $\hat{y} = S(G_e(x_s))$ for the binary segmentation case. Recall that in a multi-class segmentation setting, the loss used to train the segmenter would be the *cross-entropy* loss defined in Eq. (2). It is worth mentioning that, while minimizing the segmentation loss, gradients are also backpropagated through G_e .

Feature matching. With the previous components, our model is capable of segmenting images by learning from the source dataset \mathcal{D}_s while obtaining signal from the target images in \mathcal{D}_t through the image translation task. However, this does not ensure that the segmenter will be able to segment properly the target images x_t with the features obtained from the encoder G_e . In order to endow the network with this ability, a common approach in the domain adaptation literature is to perform a matching between the encoded features from images from both domains. This matching is commonly performed in a global way, often ignoring the distribution of classes in each domain and

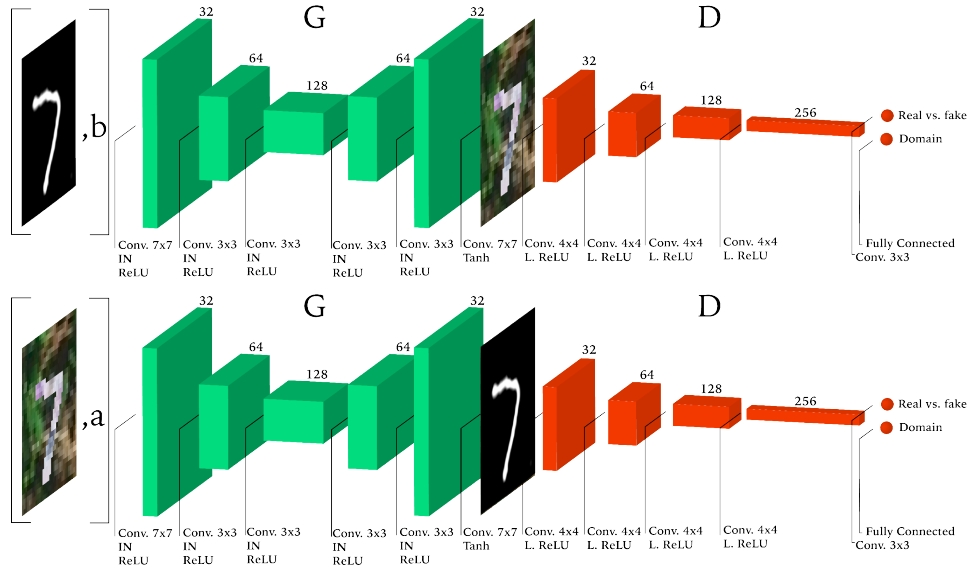


Figure 4. Architecture of the *StarGAN* [10] presented in Section III-C. *IN* indicates instance normalization and *L. ReLU* leaky rectified linear unit.

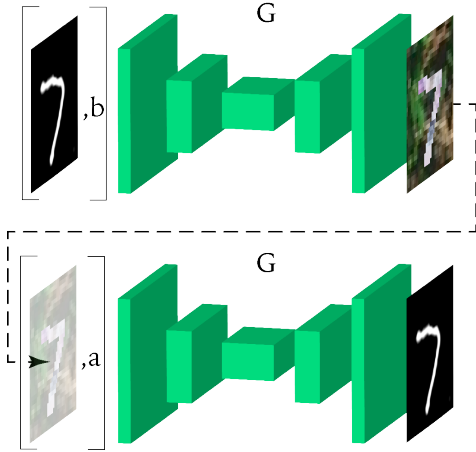


Figure 5. Cycle configuration of the modules in *StarGAN* [10] in order to compute \mathcal{L}_{cyc} for domain *a*.

thus, in the case of semantic segmentation, failing to consider the domain-specific shapes and sizes of the elements. In this thesis, we introduce a class conditional adversarial feature matching with the aim of easing the above-mentioned problem.

We define a discriminator network D_f that upsamples the features $h_s = G_e(x_s, s)$ and $h_t = G_e(x_t, s)$ to their original spatial resolution, with a single channel output. D_f is implemented by two upsampling convolution and a final convolutional layer and its role is to predict, whether each pixel in the image space belongs to the source domain or the target domain. Then, the segmentation prediction from $S(h_\bullet)$ is applied to the output of the discriminator D_f with the goal of obtaining per class probabilities of features coming from source and target domains, respectively. These per class probabilities are computed as follows:

$$\frac{\sum_p \hat{y}_\bullet^{c,p} \cdot D_f(h_\bullet)^{c,p}}{\sum_p \hat{y}_\bullet^{c,p}}, \quad (17)$$

where $D_f(h_\bullet)^{c,p}$ denotes the value of the output of $D_f(h_\bullet)$ for a certain class c and pixel p . Analogously, we denote as $\hat{y}_\bullet^{c,p}$ the value of the segmentation prediction $\hat{y}_s = S(h_s)$ for the same given

pixel and class. Note that this formulation are independent of the domain and will be used with both source and target features and segmentations. A representation of the output post-processing and the details of D_f are shown in Figure 7.

Thus, the objective of D_f is to confidently distinguish the original domain of the features per class, whereas the objective of G_e is to effectively fool D_f by trying to match the per class distribution of features. Following this objective, we define the loss to be optimized by D_f , denoted $\mathcal{L}_{dom}^{D_f}$, as:

$$\mathbb{E}_{h_t, c} \left[\frac{\sum_p \hat{y}_t^{c,p} \cdot D_f(h_t)^{c,p}}{\sum_p \hat{y}_t^{c,p}} \right] - \mathbb{E}_{h_s, c} \left[\frac{\sum_p \hat{y}_s^{c,p} \cdot D_f(h_s)^{c,p}}{\sum_p \hat{y}_s^{c,p}} \right] + \lambda_{gp} \cdot \mathbb{E}_{\bar{h}, c} \left[\left(\left\| \nabla_{\bar{h}} \left(\frac{\sum_p \bar{y}^{c,p} \cdot D_f(\bar{h})^{c,p}}{\sum_p \bar{y}^{c,p}} \right) \right\|_2 - 1 \right)^2 \right], \quad (18)$$

Note that this loss follows the *Wasserstein with gradient penalty* [19] formulation, where \bar{h} and \bar{y} are sampled uniformly along a straight line between a pair of source and target features, and segmentations, respectively. Following previous notation, λ_{gp} represents the hyperparameter controlling the relative importance of the *gradient penalty* component.

At the same time, G_e is trained by minimizing the following objective function $\mathcal{L}_{dom}^{G_e}$:

$$\mathbb{E}_{h_s, c} \left[\frac{\sum_p \hat{y}_s^{c,p} \cdot D_f(h_s)^{c,p}}{\sum_p \hat{y}_s^{c,p}} \right] - \mathbb{E}_{h_t, c} \left[\frac{\sum_p \hat{y}_t^{c,p} \cdot D_f(h_t)^{c,p}}{\sum_p \hat{y}_t^{c,p}} \right]. \quad (19)$$

Note that the gradients obtained from the segmentation predictions should neither be backpropagated through S nor G_e , as the segmentations are only used for the aggregation of the values per class.

Overall, the model is trained by minimizing the joint objectives of each module. More specifically, G and S are trained by minimizing

$$\mathcal{L}^{G, S} = \lambda_{rf} \cdot \mathcal{L}_{rf}^G + \lambda_{dom} \cdot \mathcal{L}_{dom}^G + \lambda_{cyc} \cdot \mathcal{L}_{cyc} + \lambda_{segm} \cdot \mathcal{L}_{segm} + \lambda_{dom}^f \cdot \mathcal{L}_{dom}^{G_e}, \quad (20)$$

where λ_{rf} , λ_{dom} , λ_{cyc} , λ_{segm} and λ_{dom}^f control the relative importance of the five objectives; emphasizing that gradients obtained for segmentation predictions in $\mathcal{L}_{dom}^{G_e}$ are not backpropagated.

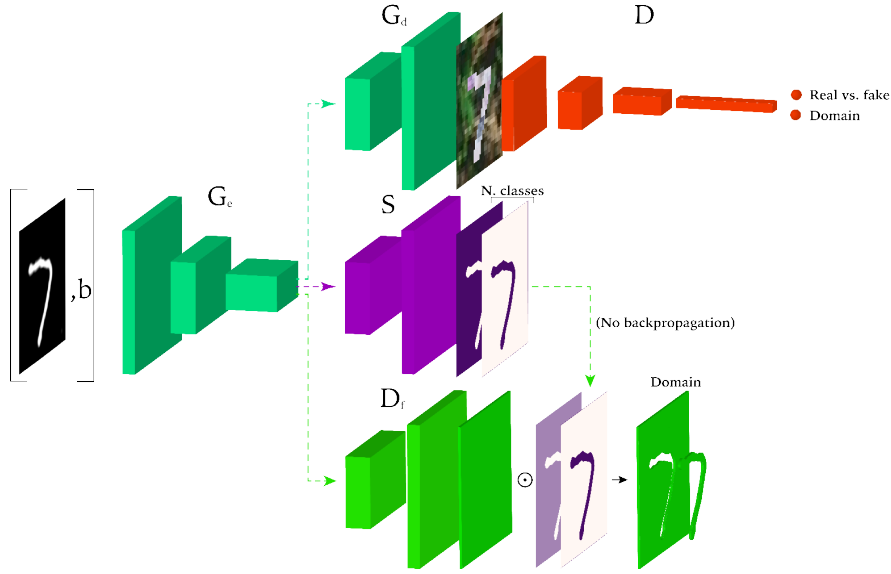


Figure 6. Complete view of the model, with the connections between the different modules: the *StarGAN* generator G (decomposed in an encoder G_e and a decoder G_d), the *StarGAN* discriminator D , the segmenter S and the feature discriminator D_f .

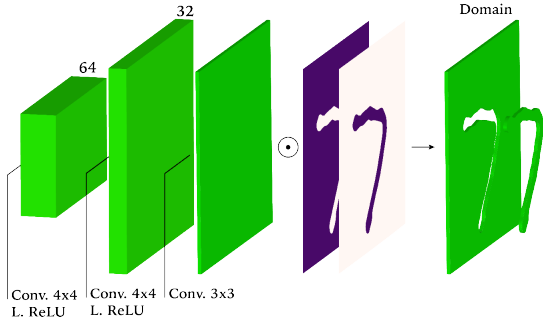


Figure 7. Structure of the feature discriminator D_f of our model, in addition to a representation of the post processing of its output. D_f upsamples the features to spatial resolution predicting whether each pixel in the image space belongs to the source domain or the target domain. Then the prediction of the segmentation masks are applied to this prediction in order to obtain the probabilities separated per class. *L. ReLU* indicates leaky rectified linear unit.

Furthermore, the discriminator D from the *StarGAN* module is trained by minimizing

$$\mathcal{L}^D = \lambda_{rf} \cdot \mathcal{L}_{rf}^D + \lambda_{dom} \cdot \mathcal{L}_{dom}^D, \quad (21)$$

where λ_{rf} and λ_{dom} are the same values used in $\mathcal{L}^{G,S}$.

Finally, we train the feature discriminator D_f by minimizing the objective

$$\mathcal{L}^{D_f} = \mathcal{L}_{dom}^{D_f}. \quad (22)$$

V. EXPERIMENTS

In this section, we will introduce and detail the experiments we have performed to measure and prove the improvements of our model with respect to some baselines and to study the relative importance of the improvement with a model ablation.

A. Datasets

The datasets used for the experiments of this thesis are selected and designed to execute a proof of concept of the contributions. Given this goal, the datasets have been prepared to simulate some features that we encounter in real scenarios, namely the differences in appearance

between different domains, as well as the per class pixel distributions. This has allowed us for a more dynamic exploration and analysis of different alternatives, which will be essential to tackle bigger and more complex datasets with a higher degree of certainty with respect to the proposed model, as future work.

MNIST. The first dataset used for the experiments is the MNIST handwritten digit dataset from LeCun et al. [31]. It contains 60,000 training samples and 10,000 test samples. Following common practice, we split the training set in 50,000 training samples and 10,000 validation samples (%16.6 of the original training set). Some samples from the dataset are shown in Figure 8a.

MNIST-M¹. We use the modification of MNIST presented in [14], which adds texture and color to the original version of MNIST, using random patches from BSDS500 dataset [1] and inverting the color in the pixels belonging to the digit. Particularly, considering I^M an MNIST image and I^B a BSDS500 random patch (both in range $[0, 1]$), the MNIST-M image is obtained as $I = |I^B - I^M|$. A few examples from the MNIST-M dataset are shown in Figure 8b. Note the varied textures and colors w.r.t. the original dataset.

The interesting characteristic about this dataset is that, when using it as a target domain in the domain adaptation problem, given a source domain as MNIST in simple black and white, it simulates the knowledge transfer from simulated environments to more complex real ones. It emulates the domain shift with colors, textures and different lighting conditions that could be found in this situations.

MNIST-thin. To obtain masks with different per class pixel distributions with respect to MNIST and MNIST-M, we generate another modification of the MNIST dataset, which is based on eroding the original digits. The process to create this new dataset is the following: we resize the MNIST images to 64×64 ; then we apply an erosion of 4 pixel radius disk and; finally, we add the skeleton of the original digit [58]. Some examples resulting from this process are shown in Figure 8c.

This modification allows us to simulate different class distributions across datasets, obtaining a %5.18 of mean foreground (digit) area per image in this dataset, compared to the original %14.64 of the MNIST and MNIST-M datasets.

¹<http://yaroslav.ganin.net/>



Figure 8. Dataset samples from MNIST, MNIST-M and MNIST-thin. Note the changes in texture, color and brightness of MNIST-M vs MNIST as well as the changes in per class pixel distributions of MNIST vs MNIST-thin.

B. Baselines

In this section, we present the baselines we have tried to overtake with our model. We train these baselines with MNIST-thin as source domain and MNIST-M as target domain, highlighting the difference in appearance and per class distribution among the two domains, and providing an insightful proof of concept. The choice of this dataset pair is motivated by the challenges it poses, namely training on simple black and white images and generalizing to a complex domain with textures and colors, as well as the challenge of dealing with different foreground/background ratios.

Single FCN. The most simple approach to overtake in the unsupervised domain adaptation problem is the one where (following the notation of Section IV) the model is trained using the source domain dataset \mathcal{D}_s , exclusively, and tested on the target domain out-of-the-box.

In our case, we train the FCN segmenter from Section III-A only using \mathcal{D}_s . In particular, we train the network on MNIST-thin and we test it on MNIST-M images. As it is a binary segmentation problem, we train the model by minimizing the *soft IOU loss* of Equation (1). The results of this setting are aimed to provide a lower bound on the performance we can achieve, given that the baseline does not include any adaptation step and due to the existing differences between the two domains.

StarGAN with segmenter. In order to exploit both \mathcal{D}_s and \mathcal{D}_t simultaneously, one could resort to training with additional auxiliary losses, e.g. for image reconstruction or image translation between domains. In this baseline, we focus on the later auxiliary task, since it has proven to be successful in the literature [37]. Note that this baseline is analogous to our model, described in Section IV, without the feature matching objectives and the class-aware discriminator D_f . The design of the model can be inferred by ignoring the D_f module in Figure 6.

In essence, the model is trained with the *StarGAN* objectives from Section III-C, together with the *soft IOU loss* from Equation (1). We train the generator G and the segmenter S by minimizing the objective

$$\mathcal{L}^{G,S} = \lambda_{rf} \cdot \mathcal{L}_{rf}^G + \lambda_{dom} \cdot \mathcal{L}_{dom}^G + \lambda_{cyc} \cdot \mathcal{L}_{cyc} + \lambda_{segm} \cdot \mathcal{L}_{segm}, \quad (23)$$

where λ_{rf} , λ_{dom} , λ_{cyc} and λ_{segm} control the relative importance of the five objectives. Following the same line, we train the discriminator D by minimizing

$$\mathcal{L}^D = \lambda_{rf} \cdot \mathcal{L}_{rf}^D + \lambda_{dom} \cdot \mathcal{L}_{dom}^D, \quad (24)$$

where λ_{rf} and λ_{dom} are the same parameters defined for $\mathcal{L}^{G,S}$.

Note that this model does not have any domain adaptation component, but should still obtain more general features than the previous single FCN baseline. For this reason, we should expect better results than the single FCN baseline, providing a more competitive lower bound.

C. Model ablation

By means of this study, we aim to assess the impact of the contributions of this thesis: the convenience of using a single image translation network with domain specific parameters and the advantage of the proposed class-conditional feature matching discriminator. Special attention will be devoted to considering alternative designs of the latter.

1) **Image-translation architectures:** In the setting of our problem, it is important to have a good image translation base model in order to obtain better features to be exploited by images from the target domain. Nevertheless, considering image translation as an auxiliary task, it is also important not to waste excessive capacity in this component (note that the decoder part of the image translator will be thrown away at test time). For this particular reason and to highlight the benefit of having domain specific encoder parameters in a single image translation network, we define this experiment to compare and contrast the results achieved by *CycleGAN* (Section III-B) w.r.t. the ones achieved by *StarGAN* (Section III-C).

Our aim is to determine whether, despite the significant reduction of number of parameters in *StarGAN* and given a similar training strategy, the resulting translations are qualitatively comparable.

2) **Class conditional discriminator:** In order to measure the benefits of conditioning the feature adversarial matching per classes (as described in Section IV), we define additional experiments using different modifications of D_f . A natural first step, is to compare our class conditional discriminator to an unconditional one. This allows us to assess the improvement achieved by our per class feature matching with respect to an often used global feature matching. Our hypothesis is that the conditioning should help in the matching due to the access to a more local discrimination. Furthermore, we define another version of a class conditioned discriminator based on the conditioning used in *StarGAN*, where the labels are given at the input. Conditioning the discriminator at input level is a natural alternative to the proposed framework, which is potentially at higher stake of making the task too easy for the discriminator. More details on these alternative discriminator designs are provided in the remainder of this section.

Unconditional discriminator. In this case, the discriminator network D_f tries to distinguish source and target domains from global features, with no class conditioning. The discriminator network is implemented by 2 downsampling convolutions, an average pooling and a fully connected classification layer. The network is depicted in Figure 9b. Note that the output of this discriminator has now spatial resolution, meaning it only outputs a single score (source vs target) for its input features. The rest of the modules remain the same as our model from Section IV.

Again, the adversarial matching is performed by means of the *Wasserstein adversarial loss with gradient penalty* [19], redefining

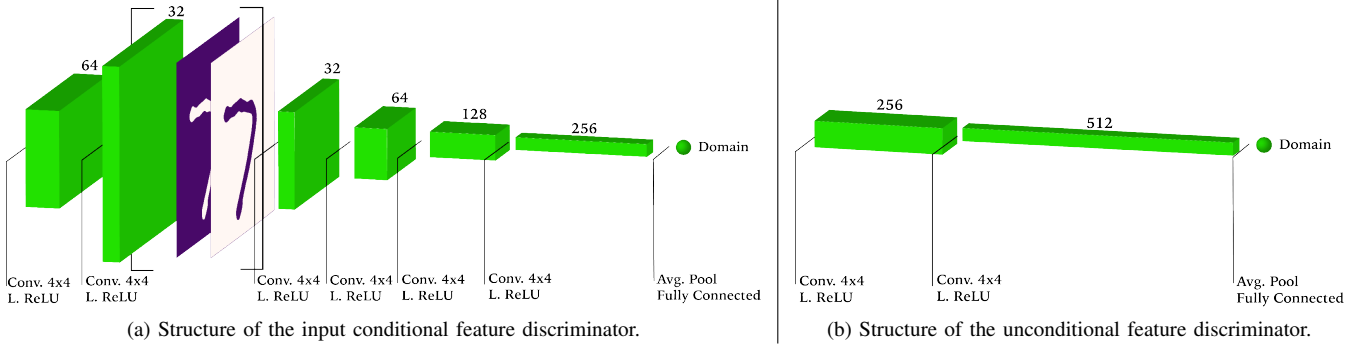


Figure 9. Alternative feature discriminators introduced in Section V-C2.

the following objectives w.r.t. the proposed model:

$$\begin{aligned} \mathcal{L}_{dom}^{D^f} &= \mathbb{E}_{h_t} [D^f(h_t)] - \mathbb{E}_{h_s} [D^f(h_s)] \\ &\quad + \lambda_{gp} \cdot \mathbb{E}_{\bar{h}} \left[\left(\left\| \nabla_{\bar{h}} D^f(\bar{h}) \right\|_2 - 1 \right)^2 \right] \\ \mathcal{L}^G &= \mathbb{E}_{h_s} [D^f(h_s)] - \mathbb{E}_{h_t} [D^f(h_t)] \end{aligned} \quad (25)$$

Input conditional discriminator. With this discriminator D_f , we aim to perform the matching of the features per class, as in our model, but providing the segmentation prediction information at input level instead. This is achieved by concatenating the segmentation prediction to the features extracted by the encoder, following the same idea as *StarGAN* with the translation labels. Note that if the encoder contains downsampling operations (like in our case), the segmentation predictions do not match the resolution of the extracted features. To overcome this situation, the network is designed with an upsampling path with 2 transposed convolutional blocks to recover the segmentation resolution. The upsampling path is applied to the features extracted by the encoder prior to concatenation with the segmentation prediction. After concatenation, the data is processed by 4 downsampling convolutional blocks followed by a fully connected layer. The architecture of the discriminator module is shown in Figure 9a. The rest of the modules remain untouched from our model.

Once again, the adversarial matching is performed by means of the *Wasserstein adversarial loss with gradient penalty* [19], redefining the following objectives from our model:

$$\begin{aligned} \mathcal{L}_{dom}^{D^f} &= \mathbb{E}_{h_t} [D^f(h_t, \hat{y}_t)] - \mathbb{E}_{h_s} [D^f(h_s, \hat{y}_s)] \\ &\quad + \lambda_{gp} \cdot \mathbb{E}_{\bar{h}} \left[\left(\left\| \nabla_{\bar{h}} D^f(\bar{h}, \bar{y}) \right\|_2 - 1 \right)^2 \right] \end{aligned} \quad (26)$$

$$\mathcal{L}_{dom}^G = \mathbb{E}_{h_s} [D^f(h_s, \hat{y}_s)] - \mathbb{E}_{h_t} [D^f(h_t, \hat{y}_t)] \quad (27)$$

VI. RESULTS

In this section, we will present and analyze the results obtained in the previously defined experiments. First, we will show qualitative results comparing *CycleGAN* and *StarGAN* as outlined in Section V-C1. Then, we will report quantitative results comparing our model to the baselines and ablated models introduced in Sections V-B and V-C, respectively.

A. Image translation results

The goal of this experiment is to assess the image translation quality of *StarGAN* vs *CycleGAN* and argue the choice of the image translation backbone in the proposed unsupervised domain adaptation model. To do so, we train both *StarGAN* and *CycleGAN* models to perform image translation between MNIST and MNIST-M

images. In order to compare both models, we provide a small set of qualitative examples in Figure 10 (additional examples are included in Appendix B). Figure 10 depicts sample translations from MNIST to MNIST-M (left) and sample translations from MNIST-M to MNIST (right) for both models. By means of these translation samples, we observe that *StarGAN* results are qualitatively comparable to those of *CycleGAN*. It is worth highlighting the variety of colors and textures in MNIST-M translations achieved by the trained *StarGAN*, as well as the rather sharp MNIST translations. *CycleGAN* exhibits slightly more uniform colors and textures, and less sharp digit translations. These details can be better perceived in the samples shown in Appendix B.

Along with the translation samples, we detail the number of parameters required by a standard implementation of both models, see Table VI-A. As shown in the table, *CycleGAN* and *StarGAN* use similar architectures for their generators and discriminators. Recall that *CycleGAN* requires two generators and two discriminators, whereas *StarGAN* only requires a single generator and a single discriminator. Therefore, *StarGAN* exhibits a reduction of roughly 50% in number of parameters when compared to *CycleGAN*. The extra parameters in the *StarGAN* generator and discriminator networks (w.r.t. the *CycleGAN* single domain networks) are caused by the domain input conditioning of the generator and the domain classification output of the discriminator. This small difference is negligible compared to the overall drop in number of parameters, when scaling the networks to higher capacities.

B. Unsupervised domain adaptation results

In this subsection, we provide both qualitative and quantitative results of the different baselines and model ablations from Section V-C, together with the results from our model. All models are trained to perform unsupervised domain adaptation from MNIST-thin to MNIST-M, i.e. only MNIST-thin segmentation masks are used for training. Recall that the goal the experiments is to assess how the proposed model performs when adapting between domains exhibiting different appearance as well as class distributions. Results are reported in terms of per class Intersection over Union (IoU) as well as mean IoU (mIoU) on both source and target domains. IoU is a common

	Generator # parameters	Discriminator # parameters
CycleGAN	$G_{ab} = 194,051$ $G_{ba} = 194,051$	$D_a = 694,241$ $D_b = 694,241$
StarGAN	197,504	694,496

Table I. Comparison of number of parameters of state-of-the-art image translation models.

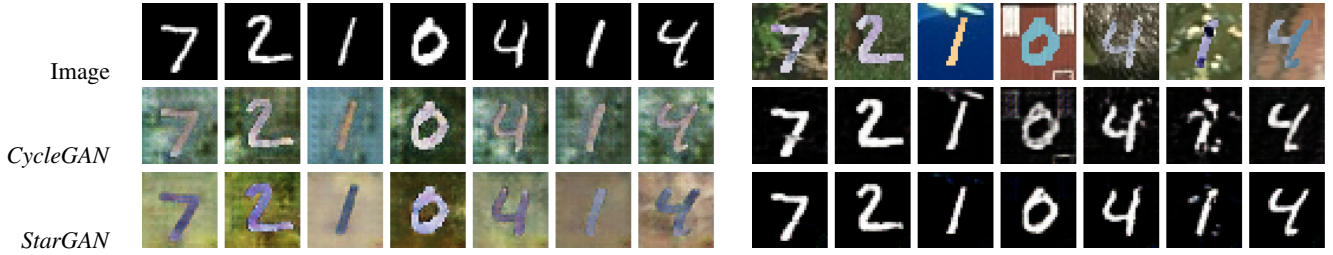


Figure 10. Example of translations obtained from *CycleGAN* and *StarGAN* between MNIST and MNIST-M datasets. Left: translations from MNIST to MNIST-M. Right: translations from MNIST-M to MNIS.

metric used in semantic segmentation and is computed as

$$\frac{TP}{TP + FN + FP}, \quad (28)$$

where TP, FN and FP denote the true positive, false negative and false positive values between a ground-truth and a prediction mask for a certain specific class. Note that given the binary segmentation masks, the maximum desired value is 1 and the minimum is 0 for each class and that the mIoU is computed as the mean over the IoU values from each class.

Table VI-B reports quantitative results and Figure 11 shows qualitative results for all methods.

We first stress the results from the *FCN segmenter* and *StarGAN with segmenter (SGAN-S)* baselines described in Section V-B. None of those baselines include any domain adaptation step, and their performance on the target domain serves solely as lower bound to the unsupervised domain adaptation models. Analogously, the performance of the *FCN segmenter* on the source domain may serve as upper bound for the source domain to the rest of the models. We also trained the same segmenter on the target domain to provide an upper bound on it; achieving a 0.998 digit IoU, a background IoU of 1 and a mean IoU of 0.999. As expected, both baselines *FCN segmenter* and *SGAN-S* exhibit the lowest results, especially when it comes to segmenting the target domain digits, outlining the need of including domain adaptation constraints to the model. However, adding an image translation backbone significantly boosts the results adding 0.127 points of digit IoU and 0.08 points of mean IoU. The differences among *FCN segmenter* and *SGAN-S* can also be visually perceived in the qualitative results shown in Figure 11, where some digits that were missed by the *FCN segmenter*, appear in the *SGAN-S* segmentations.

Second, we analyze the results of the often adopted domain adaptation strategy denoted as *SGAN-S Uncond.* (for unconditioned discriminator) in Table VI-B and Figure 11. The is *SGAN-S Uncond.* model is a natural alternative to our model, which uses a global discriminator to match features from both domains at the image translation bottleneck (see Section V-C2). We can also see *SGAN-S Uncond.* as an extension of *SGAN-S*, which includes a global feature matching as domain adaptation step. Following the reported results, we can see that the unconditional feature matching provides a notable improvement over the baselines. The target digit IoU increases almost 0.10 points and the target mIoU 0.067 points w.r.t. *SGAN-S*, which results in a 0.226 points target digit boost and 0.147 mIoU points boost over the *FCN segmenter* baseline. In the segmentation samples from Figure 11, this improvement is mainly perceived by the diminishing of false positives, leading to much cleaner segmentations (but with still room for improvement).

Third, we analyze the results of an alternative local conditional discriminator introduced in Section V-C2. We remind that in this version, the feature discriminator D_f takes the segmentation predictions as input, by concatenating them to the image translation

bottleneck features. This allows the discriminator to have information about the distributions of the classes in each domain. This model is denoted as *SGAN-S In. Cond.*. Unfortunately, this method shows a drop in performance of 0.02 points in the digit IoU over the unconditional version. We hypothesize that this performance drop is due to the differences that segmentation masks exhibit between source and target domains, i.e. source segmentations are clean and show sharp and thin digits whereas target segmentations are more noisy and are expected to show thick digits. This clear differentiation makes the task of distinguishing source and target domains trivial for the discriminator and breaks the game of the adversarial matching, showing no improvement w.r.t. a regular global adversarial feature matching.

Finally, we compare the baselines and alternative model results against our model. We denote our model as *SGAN-S + Out. Cond.*, which stands for output conditional discriminator. Recall that, in this case, the segmentation conditioning happens at the output of a pixel-wise discriminator and intends to match features per class (instead of globally). This method leads to the best results, with an improvement of 0.101 points of target digit IoU w.r.t. *SGAN-S + Uncond.*, 0.20 points of target digit IoU w.r.t. *SGAN-S* and 0.327 w.r.t. *FCN segmenter*. A similar trend can be observed for the target mIoU. Qualitatively, at the same time, the *SGAN-S + Out. Cond.* segmentations show a much better appearance, exhibiting the cleanest and sharpest segmentation predictions among all results. Moreover, the segmented digits in the target domain seem to preserve a much more realistic thickness when performing domain adaptation.

All the previously discussed evidences outline the positive impact of the contributed output conditional discriminator and support its effectiveness to properly match features representing different content, leading to improved segmentation predictions.

VII. CONCLUSIONS

In this thesis, we have tackled the unsupervised domain adaptation problem for semantic segmentation. Specifically, we have focused on the shortcomings of the global feature matching performed in most domain adaptation methods and the high numbers of parameters that state-of-the-art approaches generally need. In order to solve this problem, we have designed a new method using the information from the segmentation predictions, to perform the feature matching locally and according to the per class distributions in each domain. At the same time, we have benefited from the *StarGAN* architecture, which concatenates the domain label to the input, obtaining domain specific transformations without the need of having two completely separated domain-specific networks. With these procedures, we have been able to improve upon state-of-the-art approaches on a proof-of-concept task, while having notably less parameters to train and outlining the potential of the proposed approach. However, it is worth noting, our model still contains parameters that are only used during training, namely the image translation decoder network.

	MNIST-Thin to MNIST-M					
	Src. IoU back	Src. IoU digit	Src. mIoU	Tgt. IoU back	Tgt. IoU digit	Tgt. mIoU
<i>FCN Segmenter</i>	1	1	1	0.868	0.396	0.632
<i>SGAN-S</i>	1	1	1	0.901	0.523	0.712
<i>SGAN-S Uncond.</i>	1	1	1	0.936	0.622	0.779
<i>SGAN-S In. cond.</i>	1	1	1	0.922	0.602	0.762
<i>SGAN-S Out. cond.</i>	1	1	1	0.953	0.723	0.838

Table II. Quantitative results: test results on source (MNIST-thin) and target (MNIST-M) domain. Results are reported in terms of per class IoU (background and digit) as well as mean IoU (mIoU).



Figure 11. Segmentation samples from all models, including baselines (FCN Segmenter, SGAN-S), alternative feature discriminators (SGAN-S Uncond., SGAN-S In. Cond.), as well as our model (SGAN-S Out. Cond.).

A natural next step would be to test our method in larger scale datasets, such as SYNTHIA [44] or GTA5 [42] vs. Cityscapes [11], with the intention of analyzing the performance on more realistic environments. Additionally, we could also analyze the influence of different losses such as *identity loss* as an alternative or supplementary loss, with the goal of simplifying the model or studying loss complementarities.

In further works, it would be interesting to consider adapting our method to video processing, instead of our current image-based (frame by frame) prediction model. Exploiting the sequential nature of video, we would allow the model to output time-consistent predictions, moving closer to applications in, for example, autonomous driving systems. Another potential future research direction would be extending our methodology to one shot or zero shot learning cases where, for example, the adaptation step would need to consider the presence of unseen classes between source and target domains.

ACKNOWLEDGMENTS

First of all, I would like to thank Yaroslav Ganin for joining the project and contributing with his great ideas and experience, and Faruk Ahmed for sharing his knowledge about the adversarial networks and reviewing all our implementation. I would also like to thank Joelle Pineau for giving me the opportunity of joining the McGill RLLab and funding me during this period. Joost van de Weijer for his feedback and review of the work, and Marc Masana for his help and advice. My family, and specially my mother, for their support during this year and my roommates and classmates from Barcelona for this both awesome and hard year. And finally, I would like to thank my supervisor Adriana, first of all for offering me the opportunity of joining the AI community of Montreal and letting me work with her in this project. But, I would also like to

thank her for her hard work helping me with a really active and exceptional supervision, for taking care of my insertion in the AI research community, for teaching me so much about the appropriate research methodology (in which I had not any experience at all), for worrying about whether I was feeling comfortable in my stay and for going beyond and giving me advise for my uncertain future. Thank you.

REFERENCES

- [1] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(5):898–916, May 2011.
- [2] A. Atapour-Abarghouei and T. Breckon. Real-time monocular depth estimation using synthetic data with domain adaptation. In *2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–8. IEEE, June 2018.
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(12):2481–2495, 2017.
- [4] A. L. Bearman, O. Russakovsky, V. Ferrari, and F. Li. What’s the point: Semantic segmentation with point supervision. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VII*, pages 549–565, 2016.
- [5] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 95–104, 2017.
- [6] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 343–351, 2016.
- [7] A. Casanova, G. Cucurull, M. Drozdal, A. Romero, and Y. Bengio. On the iterative refinement of densely connected representation levels for semantic segmentation. In *2018 IEEE Conference on Computer Vision*

- and Pattern Recognition Workshops, CVPR Workshops 2018, Salt Lake City, UT, USA, June 18-22, 2018, pages 978–987, 2018.
- [8] L. Chen, J. T. Barron, G. Papandreou, K. Murphy, and A. L. Yuille. Semantic image segmentation with task-specific edge detection using cnns and a discriminatively trained domain transform. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 4545–4554, 2016.
 - [9] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
 - [10] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
 - [11] 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 *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 3213–3223, 2016.
 - [12] J. Dai, K. He, and J. Sun. Boxesup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 1635–1643, 2015.
 - [13] M. Drozdal, E. Vorontsov, G. Chartrand, S. Kadoury, and C. Pal. The importance of skip connections in biomedical image segmentation. In *Deep Learning and Data Labeling for Medical Applications - First International Workshop, LABELS 2016, and Second International Workshop, DLMIA 2016, Held in Conjunction with MICCAI 2016, Athens, Greece, October 21, 2016, Proceedings*, pages 179–187, 2016.
 - [14] Y. Ganin and V. S. Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 1180–1189, 2015.
 - [15] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, and W. Li. Deep reconstruction-classification networks for unsupervised domain adaptation. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, pages 597–613, 2016.
 - [16] R. B. Girshick. Fast R-CNN. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 1440–1448, 2015.
 - [17] A. Gonzalez-Garcia, J. van de Weijer, and Y. Bengio. Image-to-image translation for cross-domain disentanglement. *ArXiv e-prints*, May 2018.
 - [18] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014.
 - [19] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 5769–5779, 2017.
 - [20] P. Häusser, T. Frerix, A. Mordvintsev, and D. Cremers. Associative domain adaptation. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2784–2792, 2017.
 - [21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016.
 - [22] W. Hong, Z. Wang, M. Yang, and J. Yuan. Conditional generative adversarial network for structured domain adaptation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
 - [23] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2261–2269, 2017.
 - [24] N. Inoue, R. Furuta, T. Yamasaki, and K. Aizawa. Cross-domain weakly-supervised object detection through progressive domain adaptation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
 - [25] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 5967–5976, 2017.
 - [26] S. Jégou, M. Drozdal, D. Vázquez, A. Romero, and Y. Bengio. The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops, Honolulu, HI, USA, July 21-26, 2017*, pages 1175–1183, 2017.
 - [27] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II*, pages 694–711, 2016.
 - [28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
 - [29] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States.*, pages 1106–1114, 2012.
 - [30] G. Lample, N. Zeghidour, N. Usunier, A. Bordes, L. DENOYER, and M. A. Ranzato. Fader networks: manipulating images by sliding attributes. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5967–5976. Curran Associates, Inc., 2017.
 - [31] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998.
 - [32] D. Lin, J. Dai, J. Jia, K. He, and J. Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 3159–3167, 2016.
 - [33] M. Long, Z. Cao, J. Wang, and M. I. Jordan. Domain adaptation with randomized multilinear adversarial networks. *CoRR*, abs/1705.10667, 2017.
 - [34] Z. Luo, Y. Zou, J. Hoffman, and F. Li. Label efficient learning of transferable representations across domains and tasks. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 164–176, 2017.
 - [35] X. Mao, Q. Li, H. Xie, R. Y. K. Lau, Z. Wang, and S. P. Smolley. Least squares generative adversarial networks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2813–2821, 2017.
 - [36] M. Mirza and S. Osindero. Conditional generative adversarial nets. *CoRR*, abs/1411.1784, 2014.
 - [37] Z. Murez, S. Kolouri, D. J. Kriegman, R. Ramamoorthi, and K. Kim. Image to image translation for domain adaptation. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 4500–4509, 2018.
 - [38] P. O. Pinheiro. Unsupervised domain adaptation with similarity learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
 - [39] M. A. Rahman and Y. Wang. Optimizing intersection-over-union in deep neural networks for image segmentation. In *Advances in Visual Computing - 12th International Symposium, ISVC 2016, Las Vegas, NV, USA, December 12-14, 2016, Proceedings, Part I*, pages 234–244, 2016.
 - [40] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 779–788, 2016.
 - [41] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *CoRR*, abs/1804.02767, 2018.
 - [42] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *European Conference on Computer Vision (ECCV)*, volume 9906 of LNCS, pages 102–118. Springer International Publishing, 2016.
 - [43] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015*, pages 234–241, 2015.
 - [44] G. Ros, L. Sellart, J. Materzynska, D. Vázquez, and A. M. López. The SYNTHIA dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 3234–3243, 2016.

- [45] A. Rozantsev, M. Salzmann, and P. Fua. Residual parameter transfer for deep domain adaptation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [46] E. Shelhamer, J. Long, and T. Darrell. Fully convolutional networks for semantic segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(4):640–651, 2017.
- [47] R. Shu, H. H. Bui, H. Narui, and S. Ermon. A DIRT-T approach to unsupervised domain adaptation. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- [48] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [49] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [50] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1–9, 2015.
- [51] Y.-H. Tsai, W.-C. Hung, S. Schuler, K. Sohn, M.-H. Yang, and M. Chandraker. Learning to adapt structured output space for semantic segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [52] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2962–2971, 2017.
- [53] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. S. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 1349–1357, 2016.
- [54] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. S. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 1349–1357, 2016.
- [55] D. Ulyanov, A. Vedaldi, and V. S. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *CoRR*, abs/1607.08022, 2016.
- [56] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- [57] J. Zhang, Z. Ding, W. Li, and P. Ogunbona. Importance weighted adversarial nets for partial domain adaptation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [58] T. Y. Zhang and C. Y. Suen. A fast parallel algorithm for thinning digital patterns. *Commun. ACM*, 27(3):236–239, Mar. 1984.
- [59] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6230–6239, 2017.
- [60] J. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2242–2251, 2017.
- [61] J. Zhuo, S. Wang, W. Zhang, and Q. Huang. Deep unsupervised convolutional domain adaptation. In *ACM Multimedia*, 2017.

APPENDIX A MODEL IMPLEMENTATION DETAILS

In this appendix, we detail the implementation details of all networks used for the experiments of this thesis. We also provide information on the training procedure and all the necessary hyper-parameters.

Fully convolutional network for semantic segmentation. This network is explained in Section III-A and has been depicted in Figure 1. The network is built with several convolutional blocks composed of a convolution (or transposed convolution when up-sampling), a dropout layer [49], an instance normalization layer (without computation of running statistics) [55] and a rectified linear unit (ReLU) activation. More specifically, it is composed of a first convolutional block of 32 channels and 3×3 kernel size, followed by two downsampling blocks with stride 2 and 4×4 kernel size that duplicate the input channels, two upsampling blocks to recover the input resolution (dividing by 2 the number of channels) and a final 3×3 convolutional layer followed by a sigmoid non-linearity. We apply a 0.2 dropout in all convolutional blocks along the network during training.

The network is trained using an Adam optimizer [28] with an initial learning rate of 0.001 and an exponential decay of 0.995 after each epoch. The optimizer hyper-parameters are set as follows: $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Training is performed on mini-batches of size 32 for a maximum of 500 epochs with early stopping patience of 50 epochs (validation loss not improving).

CycleGAN and StarGAN. For the *CycleGAN* [60] and *StarGAN* [10] architectures presented in Sections III-B and III-C respectively, we mostly use the same implementation hyper-parameters published by the authors with slight modifications. Both models share a very similar architecture, in their generators and discriminators. The modifications that we perform result in the deletion of the bottleneck residual blocks and the reduction of the initial convolution’s number of channels to 32 in both models’ generators. These modifications aim to reduce the models’ capacity and adapt them to the datasets of interest to this thesis. In the case of *StarGAN*, we also reduce the kernel size of the generator’s convolutions to 3×3 . The models have been depicted in Figure 2 (*CycleGAN*) and Figure 4 (*StarGAN*).

Both models are trained following the optimization hyper-parameters suggested by the authors. The *CycleGAN* model is trained for 200 epochs and the best results have been picked (following a visual criteria) for the examples shown in Figure 10. In the case of *StarGAN*, the model is trained for 200.000 iterations (approximately 107 epochs).

Our model. The architecture in our model, as presented in Section IV and shown in Figure 6, is composed by several modules. The modules corresponding to the *StarGAN* backbone follow exactly the same implementation detailed previously in this appendix. The segmentation decoder S attached to G_e consists of an upsampling path and a final convolutional layer following the *FCN for segmentation* model, detailed in this appendix. The remaining module in the model is the feature discriminator D_f that we introduce in this thesis. It is depicted in detail in Figure 7 and consists of two upsampling convolutional blocks and a final 3×3 convolutional layer with a single output channel. The upsampling convolutional blocks include a 4×4 transposed convolution that upsamples the spatial resolution by 2 and outputs half of its input channels, followed by a 0.01 slope leaky ReLU. All modules have a 0.2 dropout layer in their convolutional blocks.

All the modules are trained with an Adam optimizer [28] with an initial learning rate of 0.0001 and an exponential decay of 0.995 every 1500 training iterations. The optimizer hyper-parameters are

set as follows: $\beta_1 = 0.5$ and $\beta_2 = 0.999$. The discriminators D and D_f follow a 5/1 training step ratio over the generator G and the segmenter S . The loss component importance parameters are selected as follows: $\lambda_{dom} = 1$, $\lambda_{cyc} = 10$, $\lambda_{gp} = 2$, $\lambda_{segm} = 10$ and $\lambda_{dom}^f = 1$; and finally, the whole model is trained for a maximum of 500 epochs with an early stopping patience of 50 epochs.

StarGAN with segmenter baseline. This baseline model from Section V-B follows the same architecture of our model without having a feature discriminator D_f . The training parameters are also the same as our model with the exception of the hyper-parameter $\lambda_{gp} = 10$.

Unconditional discriminator. This model presented in Section V-C2 also follows our models’ architecture, with a different discriminator D_f . This discriminator is depicted in Figure 9a and is composed of two downsampling blocks, an average pooling and a fully connected layer. The downsampling blocks consist of a 4×4 convolution with stride 2, duplicating the number of channels at its input. The convolution is followed by a leaky ReLU of slope 0.01. The module is trained with a dropout probability of 0.2.

The training setting is analogous to the one from our model, except for a λ_{gp} of 5.

Input conditioned discriminator. Analogously to the unconditional discriminator, this model also changes the discriminator D_f architecture, exclusively. The feature discriminator, in this case, is explained in Section V-C2 and depicted in Figure 9a. The architecture of this discriminator is as follows: convolutional blocks of 4×4 convolutions and leaky ReLU of 0.01 slope. In this discriminator, input features are first upsampled to the original spatial resolution by 2 upsampling convolutional blocks that double the size and halve the number of channels at each step. Then, the segmentation mask is concatenated to the upsampled features and 4 downsampling blocks are applied to reduce by a factor of 2 each axis of the spatial resolution, while duplicating the number of channels at each step (without taking into account the extra channels from the concatenation of the segmentation prediction). These features are finally processed by an average pooling, which is followed by a fully connected layer with a single output. Following previous pipelines, this discriminator is trained with a dropout probability of 0.2.

The training setting is analogous to the one followed by our model, except for a λ_{gp} of 1.

APPENDIX B ADDITIONAL EXAMPLES OF CYCLEGAN AND STARGAN TRANSLATIONS

In Figure 12 we provide additional examples of translations between MNIST and MNIST-M domains using *CycleGAN* and *StarGAN*.

APPENDIX C ADDITIONAL EXAMPLES OF SEGMENTATIONS

In Figure 13 we provide additional examples of segmentations using all models, including baselines (FCN Segmenter, *SGAN-S*), alternative feature discriminators (*SGAN-S Uncond.*, *SGAN-S In. Cond.*), as well as our model (*SGAN-S Out. Cond.*).

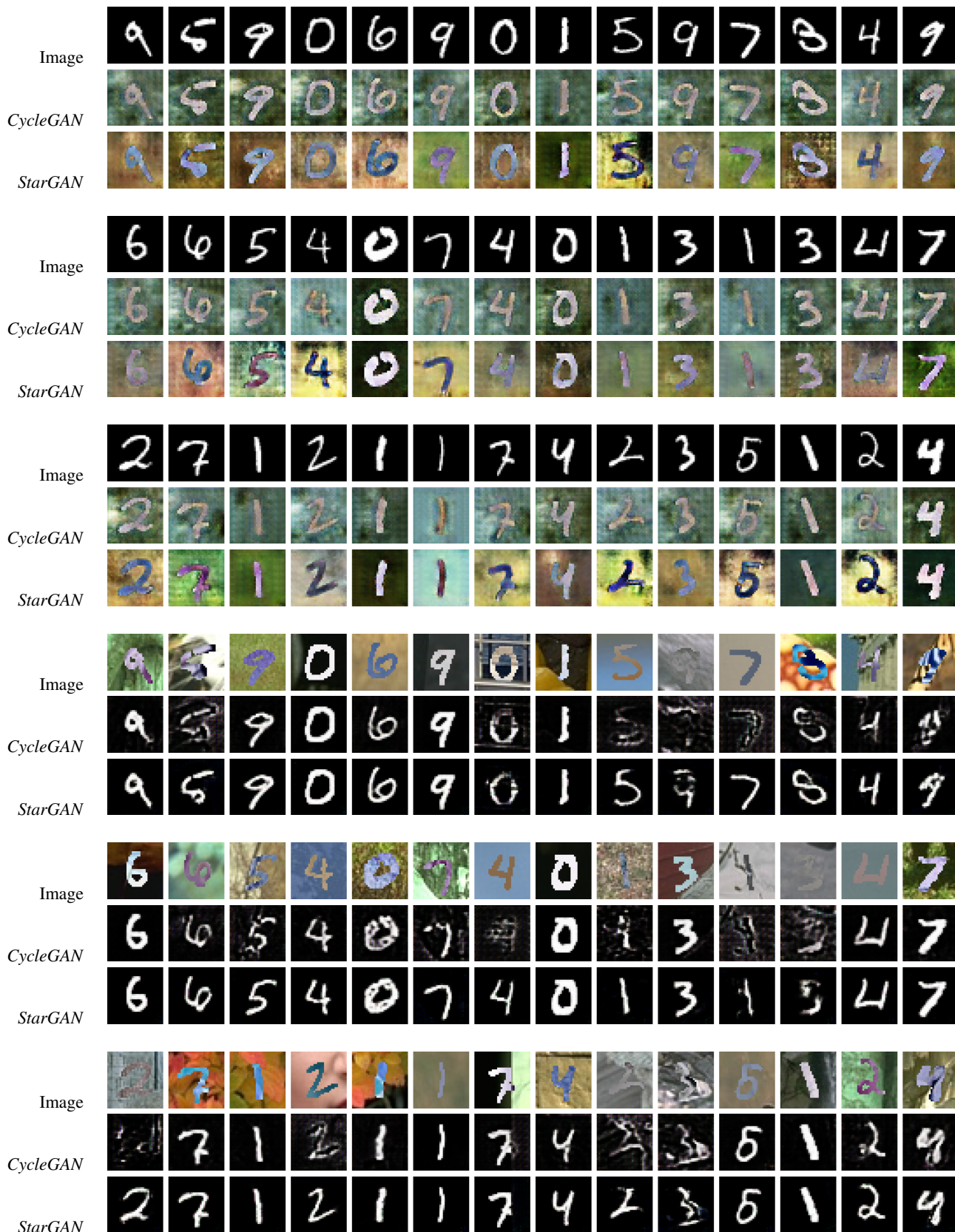


Figure 12. Additional examples of translations between MNIST and MNIST-M.



Figure 13. Additional segmentation samples from all models, including baselines (FCN Segmenter, SGAN-S), alternative feature discriminators (SGAN-S Uncond., SGAN-S In. Cond.), as well as our model (SGAN-S Out. Cond.).