

Evaluating the Impact of Image-to-Image Translation on Unsupervised Domain Adaptation Performance

Stathopoulos Dimitris
AUEB Msc In Data Science
Email: dim.stathopoulos@aueb.gr

Abstract—Unsupervised domain adaptation (UDA) aims to leverage labeled data from a source domain to improve performance in a target domain where labeled data is scarce or unavailable. In this study, we investigate the potential of UDA techniques using standard benchmarks such as MNIST, USPS, and SVHN. We evaluate the effectiveness of three models: a VGG-16 architecture augmented with Deep CORAL (Deep Correlation Alignment), a simple Convolutional Neural Network (CNN) with an integrated domain discriminator, and the Domain-Adversarial Neural Network (DANN). Despite employing these established UDA techniques, the models exhibit varying degrees of success in aligning feature representations between the source and target domains, often struggling with substantial domain shifts.

To further explore the possibilities of enhancing the adaptation process, we introduce CycleGAN for image-to-image translation, aiming to generate translated images that resemble the target domain. Our goal is to create a new source dataset that potentially bridges the domain gap. However, our experimental results indicate that the incorporation of CycleGAN for domain adaptation does not consistently yield significant improvements in target domain performance.

This study serves as a critical examination of the efficacy of combining generative models with traditional UDA techniques. While the anticipated benefits were not realized, the findings provide valuable insights into the complexities of domain adaptation and highlight the challenges faced in achieving robust and generalized solutions. Future work will continue to refine these approaches and explore alternative methods to address these persistent issues.

I. INTRODUCTION

Unsupervised domain adaptation (UDA) is a critical area in machine learning that aims to transfer knowledge from a labeled source domain to an unlabeled target domain. This task is particularly challenging when there is a significant discrepancy between the source and target domains, a scenario commonly referred to as domain shift. Various techniques have been proposed to address this issue, including feature alignment methods like Deep CORAL (Sun et al., 2016) and adversarial approaches such as the Domain-Adversarial Neural Network (DANN) (Ganin et al., 2016).

Deep CORAL reduces domain shift by aligning the second-order statistics of source and target feature distributions. In contrast, DANN employs a domain discriminator to encourage indistinguishability between source and target features, effectively creating a domain-invariant representation. Despite these advancements, achieving robust and generalized domain

adaptation remains a significant challenge, particularly for complex and highly divergent domains.

Recently, generative models like Cycle-Consistent Generative Adversarial Networks (CycleGAN) (Zhu et al., 2017) have been explored for domain adaptation. CycleGANs are capable of learning mappings between two domains without paired examples, making them an attractive option for generating synthetic data that resembles the target domain. The hypothesis is that by transforming the source domain data to appear similar to the target domain, models trained on this transformed data could perform better on the target task.

In this study, we investigate the potential of CycleGAN to enhance traditional UDA methods. We apply CycleGAN to generate target-like images from the source domain and evaluate the performance of several UDA models, including VGG-16 with Deep CORAL, a simple CNN with a domain discriminator, and DANN, on standard benchmark datasets such as MNIST, USPS, and SVHN. While previous work has shown promising results with GAN-based approaches in some contexts (Hoffman et al., 2018), our experiments reveal the nuanced and often limited benefits of this technique for UDA.

Our contributions are twofold: first, we provide an empirical evaluation of the integration of CycleGAN with traditional UDA methods; second, we offer insights into the challenges and limitations of using generative models for domain adaptation. This study aims to inform future research directions by highlighting the complexities and potential pitfalls in combining these approaches.

II. RELATED WORK

Unsupervised domain adaptation (UDA) has seen significant progress in recent years, with various techniques proposed to address the challenges of domain shift. Among the most notable approaches are those leveraging deep learning models like VGG-16 and domain alignment techniques such as Deep CORAL and adversarial methods like DANN.

A. VGG-16 w/ DeepCORAL

The VGG-16 architecture, introduced by Simonyan and Zisserman (2015), has been widely adopted in computer vision tasks due to its deep convolutional structure and ability to learn rich feature representations. To address domain adaptation, Deep CORAL (Sun Saenko, 2016) extends VGG-

16 by incorporating a correlation alignment (CORAL) loss that aligns the second-order statistics of the source and target feature distributions. This approach effectively minimizes the domain shift by ensuring that the covariance matrices of the source and target features are similar.

Deep CORAL has demonstrated promising results in various domain adaptation tasks. For instance, in the adaptation from the MNIST to USPS datasets, Deep CORAL with VGG-16 achieves competitive performance by reducing the domain discrepancy (Sun et al., 2016). Similarly, when applied to more complex datasets like Office-31, Deep CORAL has shown to significantly improve the classification accuracy by aligning feature distributions across domains (Sun et al., 2016).

B. Domain-Adversarial Neural Network (DANN)

Another prominent approach in UDA is the Domain-Adversarial Neural Network (DANN), introduced by Ganin et al. (2016). DANN employs a gradient reversal layer that allows the model to learn domain-invariant features by adversarially training a domain discriminator alongside the primary task classifier. The domain discriminator aims to distinguish between the source and target features, while the feature extractor tries to deceive the discriminator, resulting in domain-invariant feature learning.

DANN has been extensively evaluated on several benchmark datasets, demonstrating robust performance across different domain adaptation scenarios. For example, in the adaptation from SVHN to MNIST, DANN achieves high accuracy by effectively aligning the feature distributions (Ganin et al., 2016). Additionally, DANN has been successfully applied to more challenging domain adaptation tasks, such as those involving the Office-31 dataset, where it consistently outperforms traditional non-adversarial methods (Ganin et al., 2016).

Both VGG-16 with Deep CORAL and DANN represent significant advancements in UDA, offering different but complementary strategies to address the domain shift problem. While Deep CORAL focuses on aligning statistical properties of feature distributions, DANN leverages adversarial training to achieve domain invariance. These approaches have set benchmarks for acceptable performance in domain adaptation tasks, providing a foundation for further research and development.

III. METHODOLOGY

In this section, we describe the datasets, models, and experimental setup used in our study. We detail the implementation of VGG-16 with Deep CORAL, the simple CNN with a domain discriminator, and DANN. We also outline the use of CycleGAN for image-to-image translation to generate a transformed source dataset.

A. Datasets

We utilize three standard benchmark datasets for domain adaptation: MNIST, USPS, and SVHN. MNIST and USPS are both digit recognition datasets but with different styles, while SVHN consists of street view house numbers, posing a more

complex domain adaptation challenge due to its significantly different visual characteristics.

B. Models

VGG-16 with Deep CORAL: We implement the VGG-16 model modifying to accept 32x32 size images on the first Conv2D layer (algorithm 1) and extend it with the CORAL loss, which aligns the second-order statistics of the source and target domain features.

Simple CNN with Domain Discriminator: This model consists of a basic CNN for feature extraction, coupled with a domain discriminator network that classifies features as belonging to either the source or target domain. The feature extractor is trained adversarially to produce domain-invariant features.

Domain-Adversarial Neural Network (DANN): DANN incorporates a gradient reversal layer that inverts the gradient from the domain discriminator, promoting the learning of domain-invariant features by the feature extractor.

C. CycleGAN for image Translation

To further investigate the potential of enhancing UDA, we employ CycleGAN to translate images from the source domain to the target domain style. CycleGAN uses cycle-consistent adversarial networks to generate realistic target-like images from the source images without requiring paired examples.

D. Experimental Setup

We train the models using the original source datasets and the CycleGAN-transformed source datasets. We evaluate the performance of each model on the target domain using standard classification metrics. The training and evaluation protocols are consistent across all experiments to ensure a fair comparison.

1) SimpleCNN w/ Discriminator: In this work, we employ a training methodology inspired by recent advances in domain adaptation techniques using deep learning models. Specifically, we utilize a variant of the SimpleCNN architecture augmented with a domain discriminator, as outlined in Algorithm 1. This approach enables simultaneous learning of task-specific features and domain-invariant representations through adversarial training. During training, the SimpleCNN model is optimized on labeled source domain data while the domain discriminator distinguishes between source and target domain features, encouraging the model to generalize better to unseen target domain data. Our experiments demonstrate that integrating domain adversarial learning into the training process effectively enhances the model's transferability across domains, contributing to improved performance on target domain tasks.

Algorithm 1 SimpleCNN with Domain Discriminator Training Algorithm

```

1: Input: SimpleCNN model model, domain discriminator domain_model, source
   data loader src_loader, target data loader tgt_loader, loss function
   criterion, SGD optimizer opt, training parameters params, number of epochs
   epochs
2: Output: Trained SimpleCNN model with domain discriminator
3: model.train()
4: domain_model.train()
5: for epoch in range(epochs) do
6:   for (src_data, src_labels), (tgt_data, _) in zip(src_loader,
   tgt_loader) do
7:     Train feature extractor and classifier on source data
8:     opt.zero_grad()
9:     src_feat ← model.feature_extractor(src_data)
10:    src_out ← model.clf(src_feat.view(src_feat.size(0),
   -1))
11:    src_loss ← criterion(src_out, src_labels)
12:    src_loss.backward()
13:    opt.step()
14:     Train domain discriminator on source and target data
15:     opt.zero_grad()
16:     dom_labels ← torch.cat((torch.zeros(src_data.size(0)),
   torch.ones(tgt_data.size(0))).long())
17:     comb_data ← torch.cat((src_data, tgt_data))
18:     comb_feat ← model.feature_extractor(comb_data)
19:     dom_out ← dom_model(comb_feat.view(comb_feat.size(0),
   -1))
20:     dom_loss ← criterion(dom_out, dom_labels)
21:     dom_loss.backward()
22:     opt.step()
23:   end for
24: end for
   =0

```

Fig. 1. Pseudocode for training SimpleCNN with a domain discriminator.

2) *Pre-Training VGG-16*: In our approach, we pre-trained the VGG-16 model on the MNIST and SVHN datasets following a structured algorithm comprising both training and validation steps (see Algorithm 2). During each epoch, the model was trained using the training data loader with stochastic gradient descent optimization. For each batch, the loss was computed using a specified criterion, followed by backpropagation and parameter updates. Post-training, the model's performance was evaluated using a validation data loader, calculating the validation loss and accuracy without updating the model's parameters. This ensured that the VGG-16 model was effectively pre-trained, balancing both learning and generalization across the dataset.

3) *Pre-Trained VGG-16 w/ DeepCORAL*: In our study, we adopt a training strategy leveraging the CORAL (CORrelation ALignment) loss to enhance domain adaptation in deep learning models. Algorithm 3 outlines our training procedure, where we simultaneously train a neural network on source domain (SVHN) and target domain (MNIST) data. During each iteration, the model is optimized using the cross-entropy loss for classification on the source domain, while incorporating the CORAL loss to minimize domain discrepancy between feature distributions of the source and target domains. This adversarial approach helps the model learn domain-invariant representations, crucial for improving performance on target domain tasks. By integrating CORAL loss into our training framework, we aim to achieve robust and generalized deep learning models capable of handling domain shift challenges

Algorithm 2 VGG-16 Training Algorithm with Validation

```

1: Input: VGG-16 model vgg16, training data loader trainloader, validation data
   loader val_loader, loss function criterion, SGD optimizer optimizer,
   training parameters params
2: Output: Trained VGG-16 model
3: for epoch in range(params['num_epochs']) do
4:   vgg16.train()
5:   running_loss ← 0.0
6:   for images, labels in trainloader do
7:     optimizer.zero_grad()
8:     Forward pass
9:     outputs ← vgg16(images)
10:    loss ← criterion(outputs, labels)
11:    Backward pass and optimization
12:    loss.backward()
13:    optimizer.step()
14:    running_loss ← running_loss + loss.item()
15:   end for
16:   Validation step
17:   vgg16.eval()
18:   val_loss ← 0.0
19:   correct ← 0
20:   total ← 0
21:   with torch.no_grad():
22:     for images, labels in val_loader do
23:       outputs ← vgg16(images)
24:       loss ← criterion(outputs, labels)
25:       val_loss ← val_loss + loss.item()
26:       predicted ← torch.max(outputs.data, 1)
27:       total ← total + labels.size(0)
28:       correct ← correct + (predicted == labels).sum().item()
29:     end for
30:   val_loss ← val_loss / len(val_loader)
31:   val_accuracy ← 100 × correct / total
32: end for
   =0

```

Fig. 2. Pseudocode for training VGG-16 with validation step.

effectively.

Algorithm 3 VGG-16 w/ DeepCORAL Algorithm

```

1: Input: VGG-16 w/ DeepCORAL model, source data loader src_trainloader,
   target data loader tgt_trainloader, loss function criterion, Adam optimizer
   optimizer,
2: Output: Trained VGG-16 model
3: model.train()
4: result = []
5: src, tgt = list(iter(src_trainloader)),
   list(iter(tgt_trainloader))
6: train_steps = min(len(src), len(tgt))
7: for batch_idx in range(train_steps) do
8:   src_data, src_label = src[batch_idx]
9:   tgt_data, _ = tgt[batch_idx]
10:  src_data, src_label = Variable(src_data),
   Variable(src_label)
11:  tgt_data = Variable(tgt_data)
12:  optimizer.zero_grad()
13:  out1, out2 = model(src_data), model(tgt_data)
14:  clf_loss = F.cross_entropy(out1, src_label)
15:  coral_loss = CORAL(out1, out2)
16:  total_loss = _lambda * coral_loss + clf_loss
17:  total_loss.backward()
18:  optimizer.step()
19: end for =0

```

Fig. 3. Training Procedure with CORAL Loss

4) *Preprocessing Step (CycleGAN)*: CycleGAN, or Cycle-Consistent GAN, is a type of generative adversarial network for unpaired image-to-image translation. For two domains, i.e. X is MNIST and Y is SVHN, our objective is to make MNIST look like SVHN to aid unsupervised domain

adaptation. CycleGAN learns mappings $G : X \rightarrow Y$ and $F : Y \rightarrow X$. The novelty lies in trying to enforce the intuition that these mappings should be reverses of each other and that both mappings should be bijections. This is achieved through a **cycle consistency loss** that encourages $F(G(x)) \approx x$ and $G(F(y)) \approx y$. Combining this loss with the adversarial losses on X and Y yields the full objective for unpaired image-to-image translation.

For the mapping $G : X \rightarrow Y$ and its discriminator D_Y , we have the objective:

$$\mathcal{L}_{GAN}(G, D_Y, X, Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - D_Y(G(x)))] \quad (1)$$

where G tries to generate images $G(x)$ that look similar to images from domain Y , while D_Y tries to discriminate between translated samples $G(x)$ and real samples y . A similar loss is postulated for the mapping $F : Y \rightarrow X$ and its discriminator D_X .

The **Cycle Consistency Loss** reduces the space of possible mapping functions by enforcing forward and backward consistency:

$$\mathcal{L}_{\text{cyc}}(G, F) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \|F(G(x)) - x\|_1 + \mathbb{E}_{y \sim p_{\text{data}}(y)} \|G(F(y)) - y\|_1 \quad (2)$$

The full objective is:

$$\mathcal{L}_{GAN}(G, F, D_X, D_Y) = \mathcal{L}_{GAN}(G, D_Y, X, Y) + \mathcal{L}_{GAN}(F, D_X, Y, X) + \lambda \mathcal{L}_{\text{cyc}}(G, F) \quad (3)$$

where we aim to solve:

$$G^*, F^* = \arg \min_{G, F} \max_{D_X, D_Y} \mathcal{L}_{GAN}(G, F, D_X, D_Y)$$

5) *DANN Training*: In this study, we employ a comprehensive training loop for both Source Only and Domain-Adversarial Neural Network (DANN) models, as illustrated in Algorithm 4. The algorithm begins by initializing the encoder (E), classifier (C), and discriminator (D) for DANN, along with the respective source and target training loaders (S_{train} and T_{train}) and necessary parameters. For each epoch, the encoder and classifier, and discriminator in the case of DANN, are set to training mode. The algorithm calculates the current progress (p) to adjust the learning rate scheduler and the domain adaptation parameter (α). For each batch of source and target data, classification loss is computed using the source data, and domain loss is computed using both source and target data if using DANN. The total loss, comprising classification and domain losses, is backpropagated to update the model parameters. The training process includes periodic progress logging and concludes with an evaluation of the trained models on both source and target test datasets. This method ensures robust training and domain adaptation, facilitating improved generalization across domain Training Procedure with CORAL Loss

Algorithm 4 Training Loop for Source Only and DANN

```

0: Input: Encoder  $E$ , Classifier  $C$ , Discriminator  $D$  (for DANN), Source Training
  Loader  $S_{\text{train}}$ , Target Training Loader  $T_{\text{train}}$ , Parameters  $params$ 
0: Output: Trained Models  $E$ ,  $C$ , (and  $D$  for DANN)
0: Initialize loss functions:  $\text{Loss}_{\text{class}} \leftarrow \text{CrossEntropyLoss}()$ 
0: Initialize optimizer with parameters from  $E$ ,  $C$  (and  $D$  for DANN)
0: for epoch  $e$  from 1 to  $params['epochs']$  do
0:   Set  $E$  and  $C$  (and  $D$  for DANN) to training mode
0:    $start\_steps \leftarrow e \times \text{len}(S_{\text{train}})$ 
0:    $total\_steps \leftarrow params['epochs'] \times \text{len}(T_{\text{train}})$ 
0:   for batch  $(b_s, b_t)$  in  $\text{zip}(S_{\text{train}}, T_{\text{train}})$  do
0:      $(s\_images, s\_labels) \leftarrow b_s$ 
0:      $(t\_images, t\_labels) \leftarrow b_t$ 
0:      $p \leftarrow \frac{(batch\_idx + start\_steps)}{total\_steps}$ 
0:     optimizer  $\leftarrow$  optimizer_scheduler(optimizer,  $p$ )
0:     optimizer.zero_grad()
0:      $s\_features \leftarrow E(s\_images)$ 
0:     Classification Loss:
0:      $s\_class\_pred \leftarrow C(s\_features)$ 
0:      $class\_loss \leftarrow \text{Loss}_{\text{class}}(s\_class\_pred, s\_labels)$ 
0:     if Using DANN then
0:        $alpha \leftarrow \frac{2}{1 + \exp(-10 \times p)} - 1$ 
0:        $combined\_images \leftarrow \text{concat}(s\_images, t\_images)$ 
0:        $combined\_features \leftarrow E(combined\_images)$ 
0:       Domain Loss:
0:        $domain\_pred \leftarrow D(combined\_features, alpha)$ 
0:        $domain\_labels \leftarrow \text{concat}(\text{zeros\_like}(s\_labels), \text{ones\_like}(t\_labels))$ 
0:        $domain\_loss \leftarrow \text{Loss}_{\text{class}}(domain\_pred, domain\_labels)$ 
0:        $total\_loss \leftarrow class\_loss + domain\_loss$ 
0:     else
0:        $total\_loss \leftarrow class\_loss$ 
0:     end if
0:     Backpropagate:  $total\_loss.backward()$ 
0:     Update model weights: optimizer.step()
0:     if  $(batch\_idx + 1) \% 100 == 0$  then
0:       Print progress
0:     end if
0:   end for
0:   Evaluate: tester( $E$ ,  $C$ , ( $D$ ),  $S_{\text{test}}$ ,  $T_{\text{test}}$ , 'Source_only' or 'DANN')
0: end for
0: Save models and generate visualizations
=0

```

IV. EXPERIMENTAL RESULTS

We conduct extensive experiments to evaluate the effectiveness of VGG-16 with Deep CORAL, the simple CNN with domain discriminator, and DANN on the task of domain adaptation. We present quantitative results in terms of classification accuracy and provide qualitative analyses of the feature alignments achieved by each method.

A. VGG-16 Pre-Training

Pre-training of the VGG-16 model took place for MNIST and SVHN datasets separately, thus, resulting in two in total pre-trained VGG-16 classifiers. The train, validation and test sets were comprised of 84% of the original training set (from torchvision) with the remaining percentage acting as a validation set and the test set was 10000 samples, for both datasets. The pre-training results for the VGG-16 model on

VGG-16	Validation	Set Test Set
Loss	0.0105	0.0101
Accuracy	99.73%	99.78%

TABLE I
MNIST VGG-16 PRE-TRAINING RESULTS.

the MNIST dataset are summarized in Table ?? . The model achieved an impressive validation loss of 0.0105 and a test

set loss of 0.0101. In terms of accuracy, the model performed exceptionally well, with a validation accuracy of 99.73% and a test set accuracy of 99.78%. These results indicate that the VGG-16 model is highly effective at learning from the MNIST training data, generalizing well to both the validation and test sets.

VGG-16	Validation	Set Test Set
Loss	0.0624	0.0688
Accuracy	98.27%	98.16%

TABLE II
SVHN VGG-16 PRE-TRAINING RESULTS.

In contrast, the pre-training results for the VGG-16 model on the SVHN dataset, as shown in Table ??, were slightly less impressive. The validation loss was recorded at 0.0624, while the test set loss was slightly higher at 0.0688. The model achieved a validation accuracy of 98.27% and a test set accuracy of 98.16%. Although the performance on the SVHN dataset is still quite strong, it is evident that the model had a harder time generalizing compared to the MNIST dataset, as indicated by the higher loss values and lower accuracy percentages.

B. Results on UDA

Following we present the results from our UDA operations on three different models. Specifically we use USPS, MNIST and SVHN as previously explained.

The performance of different models when the source domain is MNIST and the target domain is SVHN is presented in Table III. The baseline model achieves a high source accuracy of 94.26% but suffers a significant drop in target accuracy, reaching only 19.99%. The pre-trained VGG-16-DC model shows a marked improvement with a source accuracy of 99.07% and a target accuracy of 29.11%. The DANN model also improves upon the baseline with a source accuracy of 96.48% and a target accuracy of 21.33%, indicating better generalization to the target domain compared to the baseline but not as effective as the pre-trained VGG-16-DC.

Models	Source Accuracy	Target Accuracy
Baseline	94.26%	19.99%
Pre-Tr VGG-16-DC	99.07%	29.11%
DANN	96.48%	21.33%

TABLE III
SOURCE: MNIST, TARGET: SVHN AND VGG-16 IS PRE-TRAINED ON MNIST.

When the source domain is transformed-MNIST and the target domain is SVHN, the results change significantly as shown in Table IV. The baseline model struggles in this scenario, with a source accuracy of 14.43% and a target accuracy of 11.57%. The pre-trained VGG-16-DC model performs much better, achieving a source accuracy of 84.73% and a target accuracy of 27.55%. The DANN model shows moderate performance improvements over the baseline, with a source accuracy of 15.82% and a target accuracy of 18.34%, highlighting the challenge of domain adaptation between these datasets.

Models	Source Accuracy	Target Accuracy
Baseline	14.43%	11.57%
Pre-Tr VGG-16-DC	84.73%	27.55%
DANN	15.82%	18.34%

TABLE IV
SOURCE: TRANSFORMED-MNIST, TARGET: SVHN AND VGG-16 IS PRE-TRAINED ON MNIST.

Table V presents the results when the source domain is SVHN and the target domain is MNIST, with the VGG-16 model pre-trained on SVHN. The baseline model performs reasonably well with a source accuracy of 85.10% and a target accuracy of 55.01%. The pre-trained VGG-16-DC model exhibits the best performance, with a source accuracy of 91.21% and a target accuracy of 70.77%. The DANN model, however, shows an unusual pattern with a much lower source accuracy of 30.80% but a target accuracy almost identical to the baseline at 55.00%. This suggests that while the DANN model may be more effective at adapting to the target domain, it does so at the expense of performance on the source domain.

Models	Source Accuracy	Target Accuracy
Baseline	85.10%	55.01%
Pre-Tr VGG-16-DC	91.21%	70.77%
DANN	30.80%	55.00%

TABLE V
SOURCE: SVHN, TARGET: MNIST AND VGG-16 IS PRE-TRAINED ON SVHN.

V. CONCLUSIONS

In summary, the experiments underscore the inherent challenges of domain adaptation, particularly when transitioning between datasets with distinct characteristics. The relatively limited feature space of MNIST presents a significant hurdle, as it lacks the complexity found in more diverse datasets like SVHN. This limitation was further compounded by the absence of fine-tuning, which adversely impacted the performance of the models. Notably, some baseline models, despite their simplicity, managed to outperform more specialized architectures, highlighting the potential pitfalls of over-engineering solutions without adequate dataset-specific adjustments. These findings suggest that there is substantial room for improvement in domain adaptation techniques. Future work should focus on enhancing feature extraction processes, incorporating robust fine-tuning strategies, and exploring more sophisticated models that can better generalize across domains.