

# Hallucinating Agnostic Images to Generalize Across Domains

Fabio M. Carlucci<sup>1\*</sup>   Paolo Russo<sup>2</sup>   Tatiana Tommasi<sup>3</sup>   Barbara Caputo<sup>3,4</sup>

<sup>1</sup>Huawei Noah’s Ark Lab, London   <sup>2</sup>University of Rome Sapienza, Italy

<sup>3</sup>Politecnico di Torino, Italy   <sup>4</sup>Italian Institute of Technology

fabio.maria.carlucci@huawei.com   prusso@diag.uniroma1.it  
 {tatiana.tommasi, barbara.caputo}@polito.it

## Abstract

*The ability to generalize across visual domains is crucial for the robustness of artificial recognition systems. Although many training sources may be available in real contexts, the access to even unlabeled target samples cannot be taken for granted, which makes standard unsupervised domain adaptation methods inapplicable in the wild. In this work we investigate how to exploit multiple sources by hallucinating a deep visual domain composed of images, possibly unrealistic, able to maintain categorical knowledge while discarding specific source styles. The produced agnostic images are the result of a deep architecture that applies pixel adaptation on the original source data guided by two adversarial domain classifier branches at image and feature level. Our approach is conceived to learn only from source data, but it seamlessly extends to the use of unlabeled target samples. Remarkable results for both multi-source domain adaptation and domain generalization support the power of hallucinating agnostic images in this framework.*

## 1. Introduction

Domain Adaptation (DA) is at its core the quest for principled algorithms enabling the generalization of visual recognition methods. Given at least a source domain for training, the goal is to obtain recognition results as good as those achievable on source test data on *any* other target domain, in principle belonging to a different probability distribution. While originally defined assuming to have access to annotated data from a single source domain, and to unlabeled data from a different target domain [32], there is growing interest on how to leverage over multiple sources, and for domain generalization (DG), i.e. the case when it is not possible to access target data of any sort a priori. Algorithm-wise, three strategies have been proposed,

i.e. dealing with model [8, 21], feature [25, 30], or image adaptation [31, 17]. A basic assumption for both feature and image adaptation approaches is the existence of a shared space among domains, however only feature-based methods attempt to explicitly identify it [16, 18, 3]. In the image-based approaches, the domain generic component is always silently recombined with the specific domain style to obtain images that show the same content of the target, but with source-like appearance or vice-versa [17, 31, 24]. Moreover, although these methods have shown to be effective in the single source scenario, it is questionable whether they could be extended to multi-source DA, or to DG.

With this paper we make two contributions: (1) we introduce image adaptation for DG, (2) we propose an architecture that exploits the power of layer aggregation to hallucinate samples of the latent pixel space shared among domains. We call our method Agnostic DomAin GEneralization (ADAGE). To our knowledge it is the first solution to introduce an image-level component in an end-to-end deep learning architecture for DG and that can work seamlessly also in the multi-source unsupervised DA setting.

We start by acknowledging that the notion of visual cross-domain generic information is intuitive yet ambiguous, as ground truth examples of pure semantic images without a characteristic style do not exist. Thus, while it is possible to interpret the produced samples as capturing domain agnostic knowledge, it should be clear that they are built for the network’s benefit only and we do not expect them to be pleasant to the human eye. Practically, we let the network learn what this generic information is through a mapping guided by adversarial adaptive constraints. These constraints are applied directly on the agnostic space, rather than on standard images that always contain domain-specific information.

To realize the mapping we define a dedicated convolutional structure loosely related to a previous image colorization network [6]. The new architecture has a low number of parameters which prevents overfitting and at the same time

\*This work was done while at University of Rome Sapienza, Italy

allows to comfortably accommodate two gradient reversal layers that adversarially exploit both image and feature classification across domains. As the image domain discriminator maintains the ability to evaluate the similarity of a target image to the different source domains, it is straightforward to extend the method to multi-source DA, and learn how to bias the classification loss towards the sources that are more similar to the target.

We test ADAGE in the DG and multi-source DA scenarios, comparing against recent approaches [21, 45, 43]. In all experiments, for both settings, ADAGE significantly outperforms the state of the art. An ablation study and visualizations of the agnostic domain images complete our experimental study.

## 2. Related Work

In **single source DA**, *feature adaptation* approaches aim at learning deep domain invariant representations [25, 37, 4, 5, 30, 13, 35, 14, 33]. Other methods rely on adversarial loss functions [10, 40, 34]. Besides end-to-end trained architectures also two-step adaptive networks have shown practical advantages [41, 1]. Most of work based on *image adaptation* aims at producing either target-like source images or source-like target images, but it has been recently shown that integrating both the transformation directions is highly beneficial [31, 17]. In particular [17] combines both image and feature-level adaptation. Considering that the proposed network contains two generators, three discriminators and one classifier for a single source-target domain pair, its extension to multi-source DA, and even more to DG, is not straightforward. **Multi-source DA** was initially studied from a theoretical point of view [7]. Within the context of convnet-based approaches, the vanilla solution of collecting all the source data in a single domain is already quite effective. Only very recently two methods presented multi-source deep learning approaches that improve over this baseline. The method proposed in [43] builds over [10] by replicating the adversarial domain discriminator branch for each available source. A similar multi-way adversarial strategy is used also in [45], but this work comes with a theoretical support that frees it from the need of learning the source weights.

In the **DG setting**, no access to the target data is allowed, thus the main objective is to look across multiple sources for shared factors which are either searched at *model-level* to regularize the learning process on the sources, or at *feature-level* to learn some domain-shared representation. Deep model-level strategies are presented in [26, 22, 8]. The first work proposes a weighting procedure on the source models, while the others aim at separating the source knowledge into domain-specific and domain-agnostic sub-models either with a low-rank parametrized network or through a dedicated learning architecture with a shared backbone and

source-specific aggregative modules. A meta-learning approach was recently presented in [21]. Regarding feature-based methods, [27] proposed to exploit a Siamese architecture to learn an embedding space where samples from different source domains but same labels are projected nearby, while samples from different domains and different labels are mapped far apart. Both works [12, 23] exploit deep autoencoders for DG still focusing on representation learning. New DG approaches based on *data augmentation* have shown promising results. Both [36] and [42] propose domain-guided perturbation of the input instances in the embedding space, with the second work able to generalize to new targets also when starting from a single source.

Although a two-step DG solution involving an image-adaptive process, followed by a deep classifier with feature adversarial training is always possible [29], we go beyond this naïve strategy. Differently from GAN-based methods that need a typical alternating training between image adaptation and classification, we train the whole model of ADAGE with a single optimizer while performing adversarial training by inverting the gradient originating from two domain discriminators at image and feature level.

## 3. Agnostic Domain Generalization

We assume to observe  $i = 1 \dots S$  source domains with the  $i$ th domain containing  $N_i$  labeled instances  $\{x_j^i, y_j^i\}_{j=1}^{N_i}$ , where  $x_j^i$  is the  $j$ th input image and  $y_j^i \in \{1 \dots M\}$  is the class label. In addition we also have an unlabeled target domain whose data  $\{x_j^t\}_{j=1}^{N_t}$  might (DA) or might not (DG) be provided at training time. All the source and target domains share the same label space, but their marginal distribution is different thus inducing a domain shift. The goal of ADAGE is to achieve domain generalization by hallucinating images stripped down of domain specific information, that thus can be seen as samples of a machine-created agnostic domain. We obtain this by learning to modify the images such that it becomes impossible to identify their original source domain both from their pixels and from the extracted features, while maintaining their relevant semantic information. Figure 1 shows our architecture, consisting of two main components: (1) the *Hallucinator* block, in charge of generating the agnostic images from the input samples, and (2) the *Domain Generalizer*, that performs adaptation from the new domain. The architecture is end-to-end, meaning that the two components are interconnected and trained jointly.

The **Hallucinator (H)** modifies the input images to remove their domain-specific style. To achieve this, we got inspiration from the colorization literature and define a new structure exploiting the power of layer aggregation [44]: the output of two  $3 \times 3$  convolutional layers, each followed by Relu and Batch Normalization are stacked up with the input and propagated to every subsequent layer (see Figure 2). Specifically, the produced feature build up in size resulting

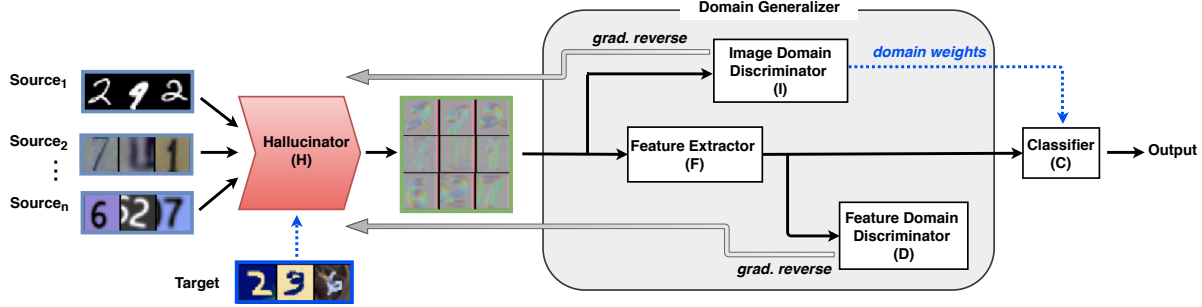


Figure 1. A schematic description of ADAGE. All samples (including target ones, in the DA setting) follow the same path in the network. The inverted gradient from  $I$  flows through  $H$  driving image modifications towards domain confusion. Similarly, the gradient from  $D$  also inverted, is backpropagated through  $F$  and  $H$  so that both the feature and the image dedicated blocks benefit from a further push towards the domain agnostic space. The classification gradient travels through the whole network, excluding  $I$  and  $D$ .

in a growing sequence of  $\{3, 8, 16, 32, 64, 128\}$  maps, after which a convolution layer brings them down to 3 channels, interpretable as RGB images. With respect to previous mapping architectures proposed within the context of depth colorization [6], our hallucinator has a significantly lower number of parameters thanks to its incrementally aggregative structure. This is crucial for generalization both because it reduces the risk of overfitting to the available sources and because leaves space for building a multi-branch following network able to impose constraints that in turn will lead to learning a stronger and more stable hallucinator in our end-to-end framework.

The **Domain Generalizer** is composed by the *Image Domain Discriminator*  $I$ , the *Feature Domain Discriminator*  $D$  and the *Feature Extractor*  $F$ . The first two impose respectively an adversarial generalization condition on the pixels and on the feature extracted from the images produced by  $H$ , while the third defines an intermediate step between the first two. Moreover, thanks to its direct connection with the *Classifier*  $C$ , it allows to maintain the basic semantic knowledge in the hallucinated images, so that, despite they lack domain style, their label can still be recognized.

The Image Domain Discriminator  $I$  receives as input the images produced by  $H$  and predicts their domain label. More in details, this module is a multi-class classifier that learns to distinguish among the  $S$  source domains in

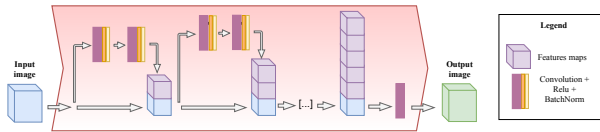


Figure 2. The Hallucinator. The output of the two multicolor blocks (Convolutional + Relu + Batch Normalization) are concatenated with the previous inputs, forming a group of images and features maps that grow along the depth of the network. The number of features increases from 3 (input data) to 256 (final aggregation step), while a last Convolutional layer squeezes the features back into 3 channels, interpretable as an RGB image.

DG, and  $S + 1$  in DA (including the target), by minimizing a simple cross-entropy loss  $\mathcal{L}_I$ . The information provided by this module is used in two ways: to adversarially guide the hallucinator  $H$  to produce images with confused domain identity, and to estimate a similarity measure between the source and the target data when available. The first task is executed through a gradient reversal layer as in [10]. The second is obtained as a byproduct of the domain classifier  $I$  by collecting the probability of every source sample in each batch to be recognized as belonging to the target.

The Feature Domain Discriminator  $D$  is analogous to  $I$  but, instead of images, it takes as input their features, performing domain classification by minimizing the cross-entropy loss  $\mathcal{L}_D$ . During backpropagation, the inverted gradient regulates the feature extraction process to confuse the domains. Finally, the Feature Extractor  $F$ , as well as the Classifier  $C$ , is a standard deep learning module. We built both of them with the same network structure used in [45] to put them on equal footing. In particular, in the DG setting the classifier learns to distinguish among the  $M$  categories of the sources by minimizing the cross-entropy loss  $\mathcal{L}_C$ , while for the DA setting it can also provide the classification probability on the target samples  $p(x^t) = C(F(H(x_t)))$  that is used to minimize the related entropy loss  $\mathcal{L}_E = p(x^t) \log(p(x^t))$ .

If we indicate with  $\theta$  the network parameters and we use subscripts to identify the different network modules, we can write the overall loss function optimized by ADAGE as:

$$\begin{aligned} \mathcal{L}(\theta_H, \theta_F, \theta_D, \theta_I, \theta_C) = & \\ & \sum_{i=1}^{S, S+1} \sum_{j=1}^{N^i} \mathcal{L}_C^{j,i}(\theta_H, \theta_F, \theta_C) + \eta \mathcal{L}_E^{j,i=S+1}(\theta_H, \theta_F, \theta_C) \\ & - \lambda \mathcal{L}_D^{j,i}(\theta_H, \theta_F, \theta_D) - \gamma \mathcal{L}_I^{j,i}(\theta_H, \theta_I). \quad (1) \end{aligned}$$

We remark that, as specified by its superscripts,  $\mathcal{L}_E^{j,i=S+1}$  is only active in the DA setting, while  $\mathcal{L}_D$  and  $\mathcal{L}_I$  in the DA case deal with an  $\{S + 1\}$ -multiclass task involving also the target together with the source domains.

As can be noted from (1), the number of meta-parameters of our approach is very limited. For  $\lambda$  we use the same rule introduced by [10] that grows the importance of the feature domain discriminator with the training epochs:  $\lambda_k = \frac{2}{1 + \exp(-10k)} - 1$ , where  $k = \frac{\text{current\_epoch}}{\text{total\_epochs}}$ . We set  $\gamma_k = 0.1\lambda_k$  so that only a small portion of the full gradient of the image domain discriminator is backpropagated: in this way we can still get useful similarity measures among the domains while progressively guiding the hallucinator to make them alike. When the image adaptation part is enough to close the domain gap, the feature discriminator loss might be abnormally high causing divergence. We easily obviate such extreme cases by maintaining a record on the initial feature discriminative loss and avoiding the loss backpropagation if it is higher than twice its initial value. Finally, the experimental evaluation indicates that ADAGE is robust to the exact choice of  $\eta$ , thus we keep it always fixed to 0.5 just for simplicity.

## 4. Experiments

We tested ADAGE<sup>1</sup> on the DG and multi-source DA scenarios. Our framework can easily switch between the two cases with a few key differences. For DG the image  $I$  and the feature  $D$  domain discriminators deal with  $S$  domains, while for DA they need to distinguish among  $S+1$  domains including the target. Moreover, in DA, the unlabeled target data trigger the classification block  $C$  to activate the entropy loss and to use the source domain weights provided by the image domain discriminator  $I$ . Specifically these weights make sure that our classifier is biased towards the sources more similar to the target.

### 4.1. Domain Generalization

**Datasets** We focus on five digits datasets and one object classification dataset. *MNIST* [20] contains 70k centered,  $28 \times 28$  pixel, grayscale images of single digit numbers on a black background. *MNIST-M* [10] is a variant where the background is substituted by a randomly extracted patch obtained from color photos of BSDS500 [2]. *USPS* [9] is a digit dataset automatically scanned from envelopes by the U.S. Postal Service containing a total of 9,298  $16 \times 16$  pixel grayscale samples; the images are centered, normalized and show a broad range of font styles. *SVHN* [28] is the challenging real-world Street View House Number dataset. It contains over 600k  $32 \times 32$  pixel color samples, while we focused on the smaller version of almost 100k cropped digits. Besides presenting a great variety of shapes and textures, images from this dataset often contain extraneous numbers in addition to the labeled, centered one. The Synthetic Digits (*SYNTH*) collection [10] consists of

500k images generated from Windows<sup>TM</sup> fonts by varying the text (that includes different one-, two-, and three-digit numbers), positioning, orientation, background and stroke colors, as well as the amount of blur. Finally, the *ETH80 object dataset* consists of 8 object classes with 10 instances for each class and 41 different views of each instance with respect to pose angles. All the images are subsampled to  $28 \times 28$  and greyscaled.

**Scenarios** We consider three experimental scenarios on digits images already presented in previous work. A first case from [45] involves **three sources** chosen in {MNIST, MNIST-M, SYNTH, SVHN}. Each dataset, with the exception of SYNTH, is used in turn as target. All the images are resized to  $28 \times 28$  pixels and subsets of 20k and 9k samples are chosen respectively from each source and from the target. A second case from [43] involves **four sources** by adding USPS to the previous dataset group, and focuses on two possible targets, SVHN and MNIST-M. Even in this case the images are resized to  $28 \times 28$  pixels, and 25/9k samples are drawn from each dataset to define the source/target sets. A third case from [12] involves **five sources** and exploits rotated variants of MNIST. Specifically we started by randomly choosing 100 images for each of the 10 classes and indicating this basic view with  $M_0$ . The versions  $\{M_{15}, M_{30}, M_{45}, M_{60}, M_{75}\}$  are obtained by rotating the images of 15 degrees in counterclock-wise direction. Note that the authors of [43] kindly shared the exact splits used for their paper, while for all the other experiments we considered multiple random selections of the samples from the datasets. For the **object classification experiment**, we followed [12] focusing on the ETH80-p setting that covers 5 domains built from equally spaced pitch-rotated views of the 8 objects. Each domain is considered in turn as the target, while the remaining ones are the sources.

**Implementation Details** For our experiments all the datasets were normalized and zero-centered. The mean and standard deviation of the target for data normalization are calculated batch-by-batch during the testing process. A standard random crop of 90 – 100% of the total image size was applied as data augmentation. The training procedure runs for 600 epochs with Adam optimizer [19]. The initial learning rate is set to  $1e^{-3}$  and step down after 80% of the training. All the experiments are repeated three times and we report the average on the obtained classification accuracy.

**Results in Table 1 (top part)** As a main baseline for the three and four sources settings we use the naïve *combine sources* strategy that consists in learning a classifier on all the source data combined together. For a fair comparison we produced these results by keeping on only the feature extractor  $F$  and the classifier  $C$ , while turning off all the adaptive blocks in the domain generalizer. We benchmark against the meta-learning method MLDG [21] using the code provided by the authors and running the experi-

<sup>1</sup>We implemented ADAGE in Pytorch, code available at <https://github.com/fmcarlucci/ADAGE/>.



Sources	SVHN	SVHN	MNIST-M	Avg.
	MNIST-M	MNIST	SYNTH	
Target	MNIST	MNIST-M	SVHN	
combine sources	98.7	62.6	69.5	76.9
DG MLDG [21]	99.1	61.2	69.7	76.7
ADAGE	99.1	<b>66.3</b>	<b>76.4</b>	<b>80.3</b>
combine sources	98.7	62.6	69.5	76.9
combine DANN [45]	92.5	65.1	77.6	78.4
MDAN [45]	97.9	68.7	81.6	82.7
ADAGE	<b>99.3</b>	<b>88.5</b>	<b>86.0</b>	<b>91.3</b>

Sources	SYNTH	SYNTH	Avg.
	MNIST	MNIST	
Target	SVHN	MNIST-M	
combine sources	73.2	61.9	67.5
DG MLDG [21]	68.0	65.6	66.8
ADAGE	<b>75.8</b>	<b>67.0</b>	<b>71.4</b>
combine sources	73.2	61.9	67.5
combine DANN [43]	68.9	71.6	70.3
DCTN [43]	77.5	70.9	74.2
ADAGE	<b>85.3</b>	<b>85.3</b>	<b>85.3</b>

Table 1. Classification accuracy results on the digits images *Left*: experiments with three sources. *Right*: experiments with four sources.

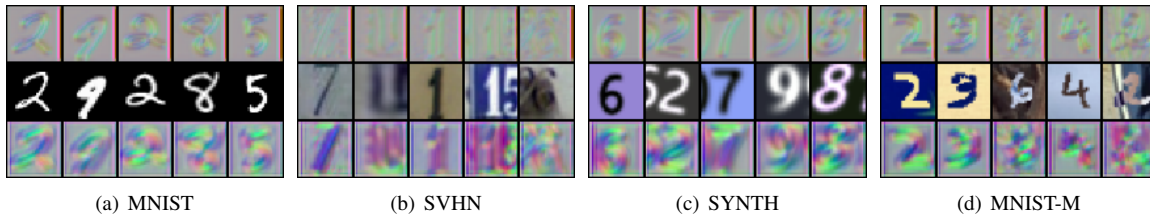


Figure 3. Examples of domain-agnostic digits generated by Hallucinator H in the three source experiments with MNIST-M as target. The top row show images produced in the DG setting by H. The central line shows the original images and in the bottom row we display images produced by H in the DA setting. **Reminder**: although we can always visualize the *domain agnostic images* to better understand the inner functioning of the network, they are not trained to be pleasant to the human eye.

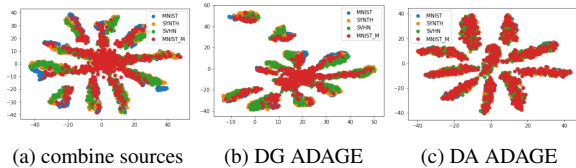


Figure 4. TSNE plots of features from the three source experiment with MNIST-M as target.

ments on our settings. The obtained results indicate that ADAGE outperforms all the reference sota baselines in DG both using three and four sources with an advantage up to 3 percentage points. Interestingly, using four sources slightly worsens the performances when SVHN is the target: our interpretation is that adding the USPS dataset increases the domain shift between the training and test domains, making the adaptation somehow more difficult.

**Results in Table 2** For the five sources experiments on rotated digit images we benchmark against two autoencoder-based DG methods D-MTAE and MMD-AAE

Target	$M_0$	$M_{15}$	$M_{30}$	$M_{45}$	$M_{60}$	$M_{75}$	Avg.
D-MTAE [12]	82.5	96.3	93.4	78.6	94.2	80.5	87.6
CCSA [27]	84.6	95.6	94.6	82.9	94.8	82.1	89.1
DG MMD-AAE [23]	83.7	96.9	95.7	85.2	95.9	81.2	89.8
CROSS-GRAD [36]	88.3	<b>98.6</b>	<b>98.0</b>	97.7	<b>97.7</b>	91.4	<b>95.3</b>
ADAGE	<b>88.8</b>	97.6	97.5	<b>97.8</b>	97.6	<b>91.9</b>	95.2

Table 2. Domain Generalization accuracy results on experiments with five MNIST-rotated sources. For compactness we only indicate the considered target.

Target	$ETH_{00}$	$ETH_{22}$	$ETH_{45}$	$ETH_{68}$	$ETH_{90}$	Avg.
combine sources	70.0	93.8	96.2	98.8	81.2	88.0
D-MTAE [12]	-	-	-	-	-	87.9
MLDG [21]	<b>70.0</b>	85.0	95.0	97.5	73.7	84.2
ADAGE	67.5	<b>95.0</b>	<b>100.0</b>	<b>100.0</b>	<b>88.8</b>	<b>90.2</b>

Table 3. Top: DG accuracy results on experiments with ETH-80 rotated sources. Bottom: real and hallucinated image examples.

respectively presented in [12] and [23], as well as against the metric-learning CCSA method [27] and the very recent CROSS-GRAD [36]. The results indicate that ADAGE outperforms three of the four competitors and has results similar to CROSS-GRAD which proposes an adaptive solution based on data augmentation that could potentially be combined with ADAGE.

**Results in Table 3** For the object classification experiments on ETH80-p, ADAGE obtains an average accuracy of 90.2% , outperforming D-MTAE [12] and MLDG [21].

## 4.2. Domain Adaptation

We extend our analysis to the multi-source DA setting considering the same three and four scenarios on digits images described in the previous section. In terms of implementation details, the only difference with respect to what already discussed above is that we now have all the un-

labeled target samples at training time, so their mean and standard deviation can be calculated at once. Moreover, for the training process we used the RmsProp optimizer [39], running for 200 epochs with initial learning rate of  $5e^{-4}$ .

**Results in Table 1 (bottom part)** We benchmark ADAGE against reference results from previous DA works. In particular for the three sources experiments the comparison is with the Multisource Domain Adversarial Network MDAN [45]. Since this method builds over the DANN algorithm [10] the result obtained with DANN applied on the combination of all the sources (combine DANN) is also reported. For the four sources experiments the main comparison is instead with the Deep Cocktail Network (DCN) [43], a recent method able to work even with partial class overlap among the sources. The results indicate that ADAGE outperforms the competing methods also in this setting with an average advantage up to 11 percentage points. As a further test we verified the obtained weights assigned by the  $I$  network component in the three source setting: when using MNIST-M as target they converge to  $\{0.5, 0.3, 0.2\}$  respectively for MNIST, SVHN, SYNTH, which sounds reasonable given the visual similarity among the domains.

While ADAGE is specifically tailored for the multi-source settings, we checked its behaviour also in the case of single source DA with access to unlabeled target data. As a proof of concept experiment, we tested ADAGE using SVHN as source and MNIST as target. With the same protocol used in our DA experiments, we achieve 95.7% accuracy, which is on par with the very recent [14] and better than several others competitive methods [17, 31, 24, 33].

### 4.3. Ablation Study and Qualitative Results

Our ablation study analyzes the effect of progressively enabling the key components of the domain generalizer alone, and in combination with the hallucinator.

**Results in Table 4** We start by evaluating the performance obtained when we do not generate the domain agnostic samples. In this case the hallucinator  $H$  is removed from the network and the original images of all sources are fed directly to the domain generalizer. In this case, since we cannot modify the original images, the only active adaptive component is  $D$  that operates on the features. Moreover the classifier can also take advantage of the entropy loss (that we indicate with  $E$ ) in the DA setting. The results indicate that feature alignment is very helpful for DA but can induce confusion in DG with results lower than those of the combine sources baseline. Another important result is obtained when only  $H$  is enabled and the features are extracted directly from the generated images with the components  $I$  and  $D$  off. In this case the network is not performing any effort to align the domains and the final accuracy is just slightly better than the combine sources baseline. This shows that the advantage of ADAGE is clearly not just due to the use of

combine sources		D	D+E	H	H+E	H+D	H+I	H+D+I	H+E+I	H+D+E	H+D+E+I	H <sub>res</sub> +D+E+I
	62.6	DG	53.0	53.0	63.2	63.2	62.2	61.4	66.3	61.4	62.2	66.3
	DA	65.9	75.1	63.2	63.9	69.9	60.8	68.8	63.9	82.4	88.5	87.6

Table 4. Ablation analysis on the experiment with three sources and target MNIST-M. We turn on and off the different parts of the model:  $H$ = Hallucinator,  $E$ = Entropy,  $D$ = Feature Domain Discriminator,  $I$ = Image Domain Discriminator. Note that  $H+D+E+I$  corresponds to our whole method ADAGE.

a deeper architecture. Keeping the hallucinator  $H$  active together with the  $D$  component produces a good advantage in accuracy but only in the DA setting ( $H + D = 69.9$ ). Here adding  $E$  provides a further advantage ( $H + D + E = 82.4$ ). Overall the entropy loss appears quite effective in the considered scenario: our intuition is that the presence of multiple sources helps reducing the risk that the entropy loss might mislead the classifier. The contribution of the image domain discriminator  $I$  is negligible by itself and this behavior can be explained considering that we backpropagate only a small part of the  $I$  gradient ( $\gamma = 0.1\lambda$ , see section 3). However its beneficial effect becomes evident in collaboration with the other network modules: passing from  $H + D + E$  to  $H + D + E + I$  implies an improvement in accuracy of at least 4 percentage points in the difficult DG setting, which shows that the adversarial guidance provided by  $I$  on  $H$  allows for an image adaptation process complementary to the feature adaptation one. Note that, since the image domain discriminator backpropagates only on the hallucinator, it is not possible to test any combination containing  $I$  but not  $H$ .

Finally we propose a benchmark against an existing residual structure previously used to transform pixels in depth image colorization [6]. When plugging in this residual version of the hallucinator ( $H_{res}$ ) we observe that the overall classification performance is slightly lower than what obtained with our original aggregative  $H$ . Besides this small variation, the most important difference is that our hallucinator has only 1/3 of the parameters of [6], thus it is faster in training and allows to avoid overfitting while mapping the source domain images into a compact agnostic space.

**Qualitative Analysis** Figure 3 shows the agnostic images generated by the hallucinator, in the three source experiment with target MNIST-M, while the bottom part of Table 3 shows examples of ETH-80 original and hallucinated images. The main effect of  $H$  is that of removing the backgrounds and enhancing the edges: this is quite clear in the DG setting for both digits and objects, while in the DA case the produced digits images appear slightly more confused. Figure 4 shows the TSNE embedding of features extracted immediately before the final classifier. In the DA setting we completely align the feature spaces of the domains, resulting in a clear per class clustering. In the DG setting the results are less clean, but the clusters are still tighter than those obtained by the combine source baseline.

## 5. Conclusions

This paper proposes the first end-to-end joint image- and feature-level adaptive solution for DG. We define a new network, named ADAGE, able to hallucinate domain agnostic images guided by two adversarial adaptive conditions at pixel and feature level. ADAGE can be seamlessly used both for DG and multi-source unsupervised DA: it achieves impressive results on several benchmarks, outperforming the current state of the art by a significant margin. We plan to extend ADAGE also to the open set multi-source DA and DG scenarios.

## Appendix

**The Hallucinator** In terms of the proposed deep learning architecture, the main technical novelty of our work is in the incremental structure of the Hallucinator module. A previous work has considered the idea of combining a sequence of consecutive layers by concatenation at the end of the network with a structure known as *hypercolumn* [15], but this approach is different from our Hallucinator: at any point in its structure we have access to all the information of the previous layers which are dynamically stacked. To further confirm the superiority of our  $H$  we extended the three sources digits experiments showing what happens when our  $H$  is substituted with a simple convolutional, residual or hypercolumn structure. The results are reported in Table 5.

### Image Domain Discriminator and Domain Weights

Note that  $I$  and  $D$  perform the same task of domain recognition, but having auxiliary losses at different levels in the network is a good strategy to better guide the learning process (e.g. see [38]). We underline that  $I$  is closer to  $H$  than  $D$ , thus  $H$  receives a cleaner signal from  $I$ . At the same time, the source-target domain similarity can be estimated better by  $I$  than  $D$ . We extended the ablation study, turning off the gradient back-propagation from  $I$  but keeping the domain similarity weights. We also moved the weights evaluation from  $I$  to  $D$ . Results (Table 6) show that adding the weights over  $H + D + E$  provides a small improvement, but not enough to outperform ADAGE.

**Target Visualization** In Figure 4 we showed two-dimensional TSNE visualization for the network features by using colors to differentiate among the domains. Here we repeat the exercise by using colors to indicate the sample class labels. The plots in Figure 5 indicate that, although the entropy loss promotes a small degree of class confusion, ADAGE clearly leads to discriminative features.

**Discussion about Related Works** While ADAGE is intuitively similar to [6] and [11], commonalities end at a high level description. The goal of [6] is to colorize raw depth images to optimally feed them to an RGB pretrained

Sources	SVHN	SVHN	MNIST-M	Avg.
	MNIST-M	MNIST	SYNTH	
Target	MNIST	MNIST-M	SYNTH	SVHN
DG $H$ Incremental (ADAGE)	99.1	66.3	<b>76.4</b>	<b>80.3</b>
$H$ Convolutional	97.9	64.0	69.5	77.1
$H$ Residual [6]	<b>99.2</b>	65.8	74.6	79.9
$H$ Hypercolumn [15]	97.7	<b>68.0</b>	70.5	78.7
DA $H$ incremental (ADAGE)	<b>99.3</b>	<b>88.5</b>	<b>86.0</b>	<b>91.3</b>
$H$ Convolutional	99.1	88.2	84.7	90.6
$H$ Residual [6]	99.2	87.6	84.1	90.3
$H$ Hypercolumn [15]	99.0	80.0	83.5	84.5

Table 5. Comparison of the incremental  $H$  architecture used in ADAGE with respect to other possible variants.

$\{\text{SVHN MNIST SYNTH}\} \rightarrow \text{MNIST-M}$	
$H + D + E$	82.4
$H + D + E + \text{weights}(I)$	84.5
$H + D + E + \text{weights}(D)$	83.3
$H + D + E + I$ (ADAGE)	<b>88.5</b>

Table 6. Extended ablation study on the DA case. In DG the weights are always off.

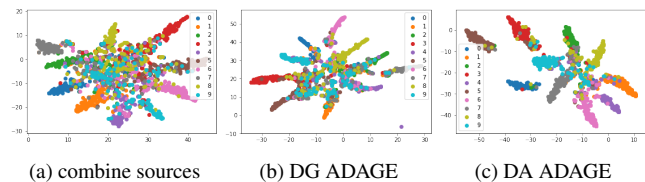


Figure 5. TSNE visualization: same setting of Fig. 5 in the main submission. Here we show only the MNIST-M target with each sample labeled according to its class.

model. Most of the network is frozen and the objective is to minimize a cross-entropy classification loss. ADAGE maps images to a domain agnostic space by minimizing two domain confusion losses besides the cross-entropy object classification loss, with the *whole* network trained end-to-end. Moreover the incremental structure of the  $H$  block is substantially different from the residual architecture used in [7]. Quantitative evidence is provided in Table 5.

The work [11] assumes that the image style is captured by the correlations between the different filter responses of any layer of the network, supposing that the style is specific for each image. ADAGE lets the network free to learn what should be considered as style over multiple images through domain confusion conditions. We make no assumptions neither on which specific part of the network captures the style, nor on how to discard it: both are automatically optimized while training across multiple domains. Indeed, by avoiding to define the style separately for any image we make it possible to deal jointly with multiple sources.

## References

- [1] G. Angeletti, B. Caputo, and T. Tommasi. Adaptive deep learning through visual domain localization. In *International Conference on Robotic Automation (ICRA)*, 2018. 2
- [2] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell. (PAMI)*, 33(5):898–916, 2011. 4
- [3] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain Separation Networks. In *Neural Information Processing Systems (NIPS)*, 2016. 1
- [4] F. M. Carlucci, L. Porzi, B. Caputo, E. Ricci, and S. Rota Bulò. Autodial: Automatic domain alignment layers. In *International Conference on Computer Vision (ICCV)*, 2017. 2
- [5] F. M. Carlucci, L. Porzi, B. Caputo, E. Ricci, and S. Rota Bulò. Just dial: domain alignment layers for unsupervised domain adaptation. In *International Conference on Image Analysis and Processing (ICIAP)*, 2017. 2
- [6] F. M. Carlucci, P. Russo, and B. Caputo. (de)2co: Deep depth colorization. *IEEE Robotics and Automation Letters*, 2018. 1, 3, 6, 7
- [7] K. Crammer, M. Kearns, and J. Wortman. Learning from multiple sources. *J. Mach. Learn. Res.*, 9:1757–1774, June 2008. 2
- [8] A. D’Innocente and B. Caputo. Domain generalization with domain-specific aggregation modules. In *German Conference on Pattern Recognition (GCPR)*, 2018. 1, 2
- [9] J. Friedman, T. Hastie, and R. Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics Springer, Berlin, 2001. 4
- [10] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *J. Mach. Learn. Res.*, 17(1):2096–2030, 2016. 2, 3, 4, 6
- [11] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *Computer Vision and Pattern Recognition (CVPR)*, 2016. 7
- [12] M. Ghifary, W. B. Kleijn, M. Zhang, and D. Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *International Conference on Computer Vision, (ICCV)*, 2015. 2, 4, 5
- [13] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, and W. Li. Deep reconstruction-classification networks for unsupervised domain adaptation. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [14] P. Haeusser, T. Frerix, A. Mordvintsev, and D. Cremers. Associative domain adaptation. In *International Conference on Computer Vision (ICCV)*, 2017. 2, 6
- [15] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *Computer Vision and Pattern Recognition (CVPR)*, 2015. 7
- [16] J. Hoffman, B. Kulis, T. Darrell, and K. Saenko. Discovering latent domains for multisource domain adaptation. In *European Conference on Computer Vision (ECCV)*, 2012. 1
- [17] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell. CyCADA: Cycle-consistent adversarial domain adaptation. In *International Conference on Machine Learning (ICML)*, 2018. 1, 2, 6
- [18] I.-H. Jhuo, D. Liu, D. T. Lee, and S.-F. Chang. Robust visual domain adaptation with low-rank reconstruction. In *Computer Vision and Pattern Recognition (CVPR)*, 2012. 1
- [19] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 4
- [20] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 4
- [21] D. Li, Y. Yang, Y. Song, and T. M. Hospedales. Learning to generalize: Meta-learning for domain generalization. In *Conference of the Association for the Advancement of Artificial Intelligence (AAAI)*, 2018. 1, 2, 4, 5
- [22] D. Li, Y. Yang, Y. Z. Song, and T. M. Hospedales. Deeper, broader and artier domain generalization. In *International Conference on Computer Vision (ICCV)*, 2017. 2
- [23] H. Li, S. Jialin Pan, S. Wang, and A. C. Kot. Domain generalization with adversarial feature learning. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 5
- [24] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *Neural Information Processing Systems (NIPS)*, 2017. 1, 6
- [25] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Deep transfer learning with joint adaptation networks. In *International Conference on Machine Learning (ICML)*, 2017. 1, 2
- [26] M. Mancini, S. R. Bulò, B. Caputo, and E. Ricci. Robust place categorization with deep domain generalization. *IEEE Robotics and Automation Letters*, 2018. 2
- [27] S. Motiian, M. Piccirilli, D. A. Adjeroh, and G. Doretto. Unified deep supervised domain adaptation and generalization. In *International Conference on Computer Vision (ICCV)*, 2017. 2, 5
- [28] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *Workshop on deep learning and unsupervised feature learning (NIPS-W)*, 2011. 4
- [29] M. M. Rahman, C. Fookes, M. Baktashmotlagh, and S. Sridharan. Multi-component image translation for deep domain generalization. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2019. 2
- [30] A. Rozantsev, M. Salzmann, and P. Fua. Beyond sharing weights for deep domain adaptation. *IEEE Trans. Pattern Anal. Mach. Intell. (PAMI)*, 2018. 1, 2
- [31] P. Russo, F. M. Carlucci, T. Tommasi, and B. Caputo. From source to target and back: symmetric bi-directional adaptive gan. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2, 6
- [32] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *European Conference on Computer Vision, (ECCV)*, 2010. 1
- [33] K. Saito, Y. Ushiku, and T. Harada. Asymmetric tri-training for unsupervised domain adaptation. In *International Conference on Machine Learning, (ICML)*, 2017. 2, 6



- [34] S. Sankaranarayanan, Y. Balaji, C. D. Castillo, and R. Chelappa. Generate to adapt: Aligning domains using generative adversarial networks. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [35] O. Sener, H. O. Song, A. Saxena, and S. Savarese. Learning transferrable representations for unsupervised domain adaptation. In *Advances in Neural Information Processing Systems (NIPS)*, 2016. 2
- [36] S. Shankar, V. Piratla, S. Chakrabarti, S. Chaudhuri, P. Jyothi, and S. Sarawagi. Generalizing across domains via cross-gradient training. In *International Conference on Learning Representations (ICLR)*, 2018. 2, 5
- [37] B. Sun, J. Feng, and K. Saenko. Return of frustratingly easy domain adaptation. In *Conference of the Association for the Advancement of Artificial Intelligence (AAAI)*, 2016. 2
- [38] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Computer Vision and Pattern Recognition (CVPR)*, 2015. 7
- [39] T. Tieleman and G. Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. In *COURSERA: Neural Networks for Machine Learning*, 2012. 6
- [40] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous deep transfer across domains and tasks. In *International Conference in Computer Vision (ICCV)*, 2015. 2
- [41] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Adversarial discriminative domain adaptation. In *Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [42] R. Volpi, H. Namkoong, O. Sener, J. C. Duchi, V. Murino, and S. Savarese. Generalizing to unseen domains via adversarial data augmentation. In *Neural Information Processing Systems (NIPS)*, 2018. 2
- [43] R. Xu, Z. Chen, W. Zuo, J. Yan, and L. Lin. Deep cocktail network: Multi-source unsupervised domain adaptation with category shift. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 4, 5, 6
- [44] F. Yu, D. Wang, E. Shelhamer, and T. Darrell. Deep layer aggregation. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [45] H. Zhao, S. Zhang, G. Wu, J. ao P. Costeira, J. M. F. Moura, and G. J. Gordon. Multiple source domain adaptation with adversarial learning. In *Workshop of the International Conference on Learning Representations (ICLR-W)*, 2018. 2, 3, 4, 5, 6