

Synthetic to Real Adaptation with Generative Correlation Alignment Networks

Xingchao Peng
Boston University
xpeng@bu.edu

Kate Saenko
Boston University
saenko@bu.edu

Abstract

Synthetic images rendered from 3D CAD models are useful for augmenting training data for object recognition algorithms. However, the generated images are non-photorealistic and do not match real image statistics. This leads to a large domain discrepancy, causing models trained on synthetic data to perform poorly on real domains. Recent work has shown the great potential of deep convolutional neural networks to generate realistic images, but has not utilized generative models to address synthetic-to-real domain adaptation. In this work, we propose a Deep Generative Correlation Alignment Network (DGCAN) to synthesize images using a novel domain adaption algorithm. DGCAN leverages a shape preserving loss and a low level statistic matching loss to minimize the domain discrepancy between synthetic and real images in deep feature space. Experimentally, we show training off-the-shelf classifiers on the newly generated data can significantly boost performance when testing on the real image domains (PASCAL VOC 2007 benchmark and Office dataset), improving upon several existing methods.

1. Introduction

Recent advances achieved by Deep Convolutional Neural Networks (DCNN) [10, 17, 36, 37, 34, 14] are unfortunately hampered by their dependence on massive amounts of training examples. Ad-hoc collection and annotation of training data for various computer vision applications is cumbersome and costly. 3D CAD simulation is a promising solution to this problem [30, 18, 35, 36, 39]. Rendering images from freely available CAD models can potentially produce an infinite number of training examples from many viewpoints and for almost any object category. Previous work [30] utilized computer graphics (CG) technique to render 2D CAD-synthetic images and consequently train deep CNN-based classifiers on them. However, their CAD-synthetic images are highly non-realistic due to the absence of natural object texture and background. More specifically, they exhibit the following problems: 1) large mismatch between foreground and background, 2) higher con-

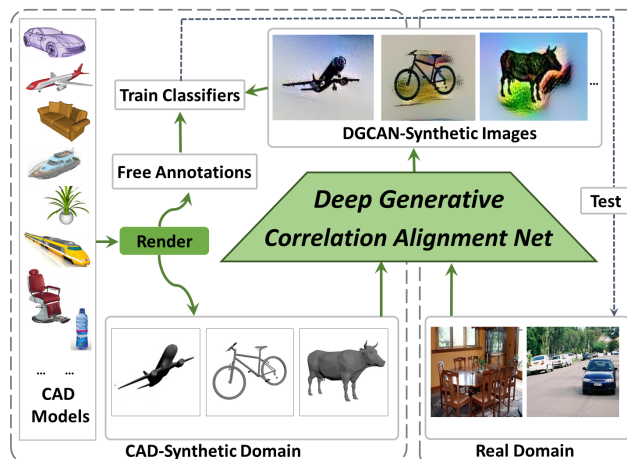


Figure 1. **Overview of our approach** We propose a Deep Generative Correlation Alignment Network (DGCAN) to bridge the domain gap between CAD-synthetic and real images in deep feature space. DGCAN can generate inexpensive annotated training data by blending the object shape from freely available 3D CAD models together with structured texture from a small amount of real background images. We train off-the-shelf classifiers on the DGCAN-synthetic images and test them on the real image domain, demonstrating a significant improvement over existing adaptation methods.

trast between the object edges and the background, 3) non-photorealistic scenery. These problems inevitably lead to a significant domain shift between CAD-synthetic and real images.

To minimize the domain shift, domain adaptation methods have been proposed to align two domains in manifold space [11, 7] or in deep feature space [20, 22, 38]. These algorithms bridge the domain gap between real-image domains. However, the domain shift between synthetic and real domains is much larger than that between two real-image domains. CAD-to-real adaptation methods [26, 5] have been proposed but they only align the viewpoint of specific indoor categories and cannot be directly applied to recognition systems in the wild.

Our main idea is to incorporate domain adaptation algorithm into generative networks. Generative neural networks have recently been proposed to create novel im-

agery that shares common properties with some given images, such as *content* and *style* [9], similarity in feature space [19, 12, 42, 24, 27], etc. However, these approaches have several limitations for use in domain adaptation. For example, *Generative Adversarial Nets* (GANs) [12] and *style transfer* approaches [9, 8] can generate images but are not designed for domain adaptation. *Coupled GANs* [19] only handle domain shifts between small images (28×28 pixel resolution). *Conditional GANs* [15] can learn image-to-image translation but need paired training data that are costly to obtain, *i.e.* CAD models and corresponding natural images.

To overcome the limitations of the aforementioned approaches, we propose a *Deep Generative Correlation Alignment Network* (*DGCAN*) to bridge the domain discrepancy between CAD-synthetic and real images. Our work is primarily motivated by [30, 9, 38]. As shown in Figure 1, we generate novel images by matching the convolutional layer features with those of a *content* CAD-synthetic image and the feature statistics of a real image containing a background scene. Unlike *neural style* [9], the goal is not to create an artistic effect but rather to adapt the CAD-synthetic data to match the statistics of real images and thus improve generalization. To this end, we employ the correlation alignment (CORAL) loss [38] for adaptation. However, instead of learning to align features, we generate images whose feature correlations match the target real-image domain.

Our synthesized results reveal that *DGCAN* can satisfactorily blend the contour of specific objects (from CAD-synthetic images) with natural textures from real images. Although the generated images are not fully photorealistic, they appear to have more natural statistics to the deep network, improving its performance. Extensive experiments on the PASCAL and Office dataset show that our approach yields a significant performance boost compared to the previous state-of-the-art methods [30, 37, 38, 7, 11, 20, 22, 9].

The contributions of this paper can be summarized as follows.

- We propose *Deep Generative Correlation Alignment Network* (*DGCAN*) to synthesize CAD objects contour from the CAD-synthetic domain with natural textures from the real image domain.
- We explore the effect of applying the content and *CORAL* losses to different layers and determine the optimal configuration to generate the most promising stimuli.
- We empirically show the effectiveness of our model over several state-of-the-art methods by testing on real image datasets.

2. Related Work

CAD Simulation CAD simulation has been extensively used by researchers since the early days of computer vision [28]. 3D CAD models have been utilized to generate stationary synthetic images with variable object poses, textures, and backgrounds [30]. Recent usage of CAD simulation has been extended to multiple vision tasks, e.g. object detection [30, 26], pose estimation [18, 35, 36, 39], robotic simulation [40], semantic segmentation [31]. However, for many tasks, CAD-synthetic images are too low-quality due to the absence of realistic backgrounds and texture. To mitigate this drawback, [30] proposes to directly add auxiliary texture and background to the rendered results, with the help of commercial software (e.g. AutoDesk 3ds MAX¹). However, this method introduces new problems, such as unnatural positioning of objects (e.g. floating car above the road), high contrast between object boundaries and background, etc. Our approach tackles these problems by synthesizing novel imagery with *DGCAN* and can generate images with natural feature statistics.

DCNN Image Synthesis Deep convolutional neural networks learn distributed, invariant and nonlinear feature representations from large-scale image repositories [1]. *Generative Adversarial Networks* (GANs) [12] and their variations [27, 29] aim to synthesize images that are indistinguishable from the distribution of images in their training set. However, training GANs is difficult and often leads to oscillatory behavior. *Style transfer* [9] synthesizes novel stimuli by aligning the *conv* layer features and *Gram Matrices* of the features. In this way, the synthesized image simultaneously preserves the arrangement of a content image (often a normal photograph) and the colours and subtle local structures of a style image (often an artist’s work). Our approach is inspired by *style transfer* [9], but is geared towards adapting a set of CAD-synthetic images to real image domain with a domain adaptation loss.

Domain Adaptation Domain shift results in a significant performance degradation when recognition systems are trained on one domain (source) and tested on another (target). Shallow domain adaptation algorithms aim to bridge the two feature distributions via mappings learned either by minimizing a distribution distance metric [2, 37], or by projecting the feature distributions to a common low-dimensional manifold [13, 11, 21]. Deep domain adaptation methods address the domain shift by adding one or multiple adaptation layers and losses [41, 40, 20, 38], or use an adversarial network to match the source distribution to target [40, 19]. All of the aforementioned methods follow the paradigm of aligning the source domain and target domain in feature space. In contrast, we take a generative approach

¹<http://www.autodesk.com/store/products/3ds-max>

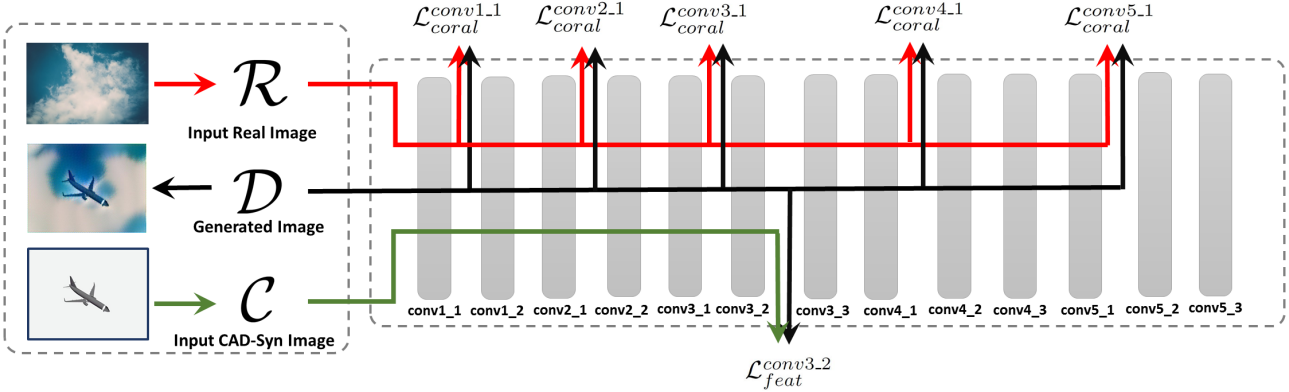


Figure 2. **Illustration of DGCAN.** Our model, the Deep Generative Correlation Alignment Network (DGCAN), takes CAD-synthetic images \mathcal{C} and real background images \mathcal{R} as input and generates novel images \mathcal{D} that contain the same object shapes as the CAD images but have more natural feature distributions by transferring the texture from the background image. The network structure is based on VGG-16 [34], which comprises 13 convolutional layers, divided in five groups. The image is generated by applying the ℓ^2 loss $\mathcal{L}_{feat}^{\mathcal{X}_f}$ and CORAL loss $\mathcal{L}_{coral}^{\mathcal{X}_c}$ to the following layers: $\mathcal{X}_f = \{\text{conv}3_2\}$, $\mathcal{X}_c = \{\text{conv}1_1, \text{conv}2_1, \text{conv}3_1, \text{conv}4_1, \text{conv}5_1\}$.

to combine the statistics of target domain images with the content of source domain images. Recent proposed generative models [3, 19] adapt two domains by adversarial losses. However, these methods only generate small images. Our model can generate large images with arbitrary resolution.

3. Approach

Suppose we are given n_s labeled source-domain CAD-synthetic (image, label) pairs $\mathcal{I}_s = \{\mathcal{C}_i, \mathcal{Y}_i\}_{i=1}^{n_s}$, and n_t target-domain real images $\mathcal{I}_t = \{\mathcal{R}_i\}_{i=1}^{n_t}$. We assume that the target domain is unlabeled, so object classifiers can only be trained on \mathcal{I}_s . However, their performance will degrade when testing on \mathcal{I}_t due to the domain discrepancy. Our aim is to synthesize a labeled intermediate dataset $\mathcal{I} = \{\mathcal{D}_i, \mathcal{Y}_i\}_{i=1}^n$, such that each $\mathcal{D}_i \in \mathcal{I}$ contains a similar object shape and contour with $\mathcal{C}_i \in \mathcal{I}_s$ and similar local pattern, color, and subtle structure (“style” as illustrated in [9]) with some random $\mathcal{R} \in \mathcal{I}_t$.

To generate \mathcal{D} from \mathcal{C} and \mathcal{R} , the most straightforward method is to average the two images. Traditional computer vision blending approaches, such as half-half alpha blending or pyramid blending lead to image artifacts that contribute to the domain shift. Previous CG-based method [30] applied real image background and texture to CAD models, leading to the problems illustrated in Section 1. Instead, we propose to align the generated \mathcal{D} to \mathcal{C} and \mathcal{R} in the DCNN feature space, as shown in Figure 2. Analogously to [9], our model synthesizes an image \mathcal{D} from \mathcal{C} with $\mathcal{D} \sim p(\mathcal{D}|\mathcal{C}, \mathcal{R})$. The generation is guided by two losses, one to ensure the object contour stays the same, and the other to ensure the image has similar low-level statistics with real images.

3.1. Deep Convolutional Neural Network

We base our approach on the VGG-16 [34] network which consists of 13 convolutional layers ($\text{conv}1_1$ - $\text{conv}5_3$), 3 fully connected layers ($\text{fc}6$ - $\text{fc}8$) and 5 pooling layers ($\text{pool}1$ - $\text{pool}5$). The convolutional layers consist of a set of learnable kernels. Each kernel is convolved with the input volume to compute hidden activations during the forward pass and its parameters are updated through a back-propagation pass. We denote $\mathcal{H}^l(\cdot)$ as DCNN’s l^{th} layer’s representation matrix, $\mathcal{H}_i^l(\cdot)$ as the i^{th} dimension of $\mathcal{H}^l(\cdot)$ and $\mathcal{H}_{ij}^l(\cdot)$ as j^{th} value of $\mathcal{H}_i^l(\cdot)$.

3.2. Shape Preserving loss

To preserve the shape information of CAD-synthetic images, we propose to use the ℓ^2 loss in feature space as follows

$$\mathcal{L}_{feat}^{\mathcal{X}_f} = \sum_{l \in \mathcal{X}_f} \left(\frac{\omega_f^l}{2\alpha^l} \sum_i \|\mathcal{H}_i^l(\mathcal{D}) - \mathcal{H}_i^l(\mathcal{C})\|_2^2 \right). \quad (1)$$

where $\mathcal{D} \in \mathcal{I}$, $\mathcal{C} \in \mathcal{I}_s$; ω_f^l is the loss weight of l^{th} layer in DCNN feature space; \mathcal{X}_f is the collection of convolutional layers which the ℓ^2 loss is applied to; $\alpha^l = N^l F^l$, where N^l is the channel number of l^{th} layer’s feature, and F^l is the length of feature in a single channel.

The derivative of this loss with respect to the activations in a particular layer l can be computed by:

$$\frac{\partial \mathcal{L}_{feat}^{\mathcal{X}_f}}{\partial \mathcal{H}_{ij}^l(\mathcal{D})} = \frac{\omega_f^l}{\alpha^l} (\mathcal{H}_{ij}^l(\mathcal{D}) - \mathcal{H}_{ij}^l(\mathcal{C})) \quad (2)$$

This gradient can be back-propagated to update the pixels while synthesizing \mathcal{D} .

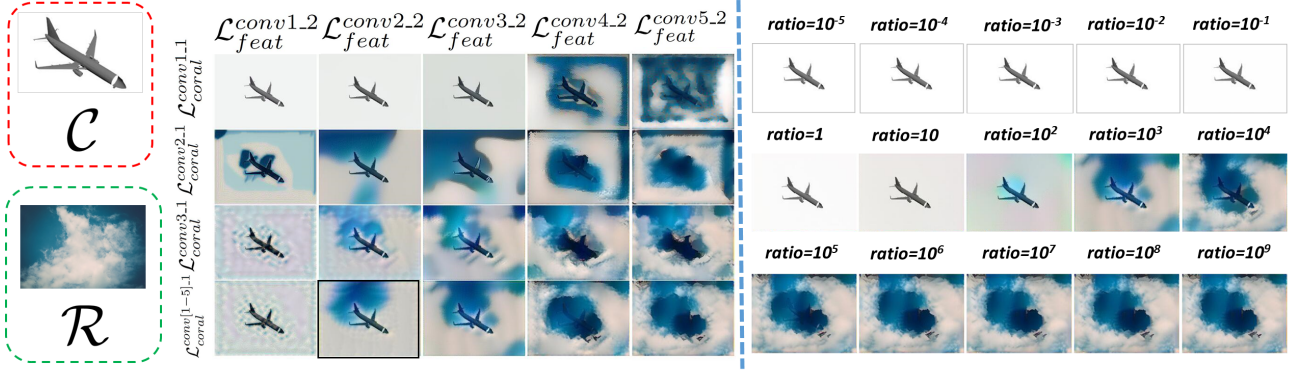


Figure 3. **Illustration of our synthesized results.** We leverage *DGCAN* to synthesize novel images based on two inputs, i.e. source domain CAD-synthetic image $\mathcal{C} \in \mathcal{I}_s$ and target domain real background image $\mathcal{R} \in \mathcal{I}_t$. **(1).** We exhaustively apply \mathcal{L}_{feat} and \mathcal{L}_{coral} to different *conv* layers to find the best configuration. The results (**left plot**) demonstrate that *DGCAN* can generate more distinct object contours when applying \mathcal{L}_{feat} to lower *conv* layers and can synthesize more structured style texture when applying the \mathcal{L}_{coral} to higher *conv* layers. ($\mathcal{L}_{coral}^{[1-5].1}$ means applying \mathcal{L}_{coral} to *conv1.1*, *conv2.1*, *conv3.1*, *conv4.1*, *conv5.1* simultaneously) **(2).** We vary the trade-off parameter λ in equation 6 from $10^{-5} \sim 10^9$ to learn the optimal value for λ . The results (**right plot**) show that the shape contour dominates the background texture when λ is small.

3.3. Naturalness loss

Networks trained on CAD images will not work well on input real images because of the mismatch in low-level image statistics such as textures, edge contrast, color, etc. To align the low-level texture statistics of the generated images to the real image domain, we propose to employ the CORAL loss. Correlation Alignment (CORAL) was first devised by [37] to match the second-order statistics of feature distributions for domain adaptation. It is derived by minimizing the domain discrepancy with squared *Frobenius* norm $\min \|Cov_S - Cov_T\|_F^2$, where Cov_S, Cov_T are the covariance matrices of feature vectors from source domain and target domain, respectively. This problem is equivalent to solving $A^* = \operatorname{argmin} \|A^T Cov_S A - Cov_T\|_F^2$.

Inspired by [37], we define the CORAL loss $\mathcal{L}_{coral}^{\mathcal{X}_c}$ as

$$\mathcal{L}_{coral}^{\mathcal{X}_c} = \sum_{l \in \mathcal{X}_c} \left(\frac{\omega_c^l}{4\alpha^l} \|Cov(\mathcal{H}^l(\mathcal{D})) - Cov(\mathcal{H}^l(\mathcal{R}))\|_F^2 \right) \quad (3)$$

where $\mathcal{D} \in \mathcal{I}$, $\mathcal{R} \in \mathcal{I}_t$; ω_c^l is the *CORAL* loss weight of l^{th} layer; \mathcal{X}_c is the collection of convolutional layers that the *CORAL* loss is applied to; $Cov(\cdot)$ is the covariance matrix of l^{th} layer's activation; $\|\cdot\|_F$ denotes the *Frobenius* distance.

Analogous to [38], the covariance matrices are given by:

$$Cov(\mathcal{H}^l(\mathcal{M})) = \frac{1}{N^l} \{ \mathcal{H}^l(\mathcal{M})^\top \mathcal{H}^l(\mathcal{M}) - \frac{1}{N^l} [(\mathbf{1}^\top \mathcal{H}^l(\mathcal{M}))^\top (\mathbf{1}^\top \mathcal{H}^l(\mathcal{M}))] \} \quad (4)$$

where $\mathcal{M} \in \{\mathcal{D}, \mathcal{R}\}$, $\mathbf{1}$ is a column all-one vector, and N^l is the number of feature channels in l^{th} layer.

The derivative of the *CORAL* loss with respect to a particular layer l can be calculated with chain rule:

$$\frac{\partial \mathcal{L}_{coral}^{\mathcal{X}_c}}{\partial \mathcal{H}_{ij}^l(\mathcal{D})} = \frac{\omega_c^l}{N^l \alpha^{l^2}} \{ [\mathcal{H}^l(\mathcal{D})^\top - \frac{1}{N^l} (\mathbf{1}^\top \mathcal{H}^l(\mathcal{D})) \mathbf{1}^\top]^\top \cdot (Cov^l(\mathcal{D}) - Cov^l(\mathcal{R})) \}_{ij} \quad (5)$$

Our final method combines the loss functions defined by equation 1 and equation 3. We start from an image $\mathcal{D} \in \mathcal{I}_s$ and pre-process it by adding a random perturbation ϵ , where $\epsilon \sim \mathcal{N}(0, \Sigma)$. We then feed the image forward through *DGCAN* and compute the ℓ^2 loss with respect to \mathcal{D} and *CORAL* loss with respect to \mathcal{R} . The back-propagated gradient thus guides the image synthesis process. Hence, the synthesized image is the output of the function:

$$\mathcal{D}^* = \operatorname{argmin}_{\mathcal{D} \in \mathcal{I}} (\mathcal{L}_{feat}^{\mathcal{X}_f} + \lambda \mathcal{L}_{coral}^{\mathcal{X}_c} | \mathcal{C}, \mathcal{R}, \mathcal{X}_f, \mathcal{X}_c, \lambda, \epsilon) \quad (6)$$

where $\mathcal{L}_{feat}^{\mathcal{X}_f} + \lambda \mathcal{L}_{coral}^{\mathcal{X}_c}$ denotes the total loss of *DGCAN*, and λ denotes the trade-off weight between ℓ^2 loss and the *CORAL* loss. The hyperparameter λ is set through cross validation.

4. Experiments

Our experiments include two parts. First, we apply *DGCAN* to the CAD-synthetic dataset provided by [30] to synthesize adapted *DGCAN*-synthetic images. Second, we train off-the-shelf classifiers on the *DGCAN*-synthetic images and test on the PASCAL 2007 [6] and Office [32] datasets. We implement our model with the *Caffe* [16]

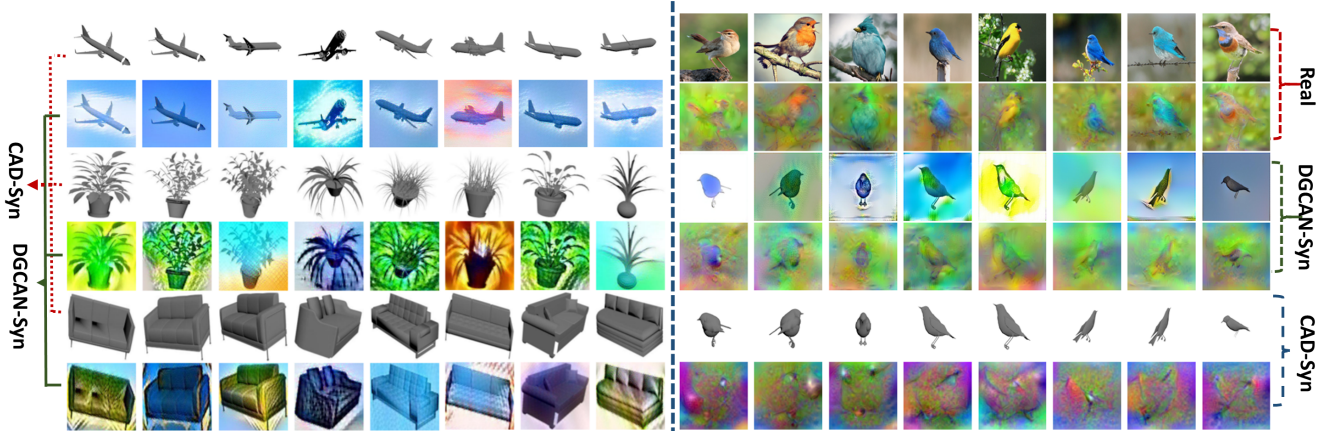


Figure 4. **DGCAN-synthetic examples and visualizations.** (1) The left plot shows randomly selected DGCAN-synthetic examples (\mathcal{D}) and their corresponding CAD-synthetic images (\mathcal{C}). The rendered results demonstrate that *DGCAN* can synthesize novel images with clear object contours and photo-realistic textures. (2) The right plot illustrates the reconstruction results generated by using the tools provided by [25]. The reconstructions reveal that our DGCAN-synthetic images share more similarities with real images from the DCNN’s perspective. The (uniform gray-scale) CAD-synthetic images only provides edge information. Thus, the pixels in the reconstructed images are dominated by the rich color and texture information encoded in the DCNN’s parameters. (Best viewed in color!)

framework. Datasets (both CAD-synthetic and DGCAN-synthetic), code and experimental configurations will be made available publicly.

4.1. Generating Adapted Images

As shown in Figure 2, while generating the DGCAN-synthetic dataset, we set CAD-synthetic images as \mathcal{C} and real images downloaded from the Google image search engine as \mathcal{R} .

CAD-Synthetic Dataset The CAD-synthetic dataset in [30] was rendered from 3D CAD models for zero-shot or few-shot learning tasks. The dataset contains 6 subsets with different configurations (i.e. RR-RR, W-RR, W-UG, RR-UG, RG-UG, RG-RR). The process of rendering the dataset (we refer the reader to [30] for more details) can be summarized as follows: (1) collecting 3D-CAD models from large-scale on-line repositories (Google Sketchup, Stanford 3D ShapeNet²), (2) selecting image cues (background, texture, pose, etc.), (3) rendering synthetic images with AutoDesk 3ds Max. In our experiments, we only adopt images with white background because other subsets suffer from the issues described in Section 1.

Parameter tuning To determine the optimal configuration for \mathcal{X}_f , \mathcal{X}_c and λ , we exhaustively apply \mathcal{L}_{feat} , \mathcal{L}_{coral} to different *conv* layers and vary λ from $10^{-5} \sim 10^9$ on a small validation dataset.

Results and Analysis A representative subset of rendered results with different setting for \mathcal{X}_f , \mathcal{X}_c and λ are shown in Figure 3. The left plot shows the effect of different configurations of \mathcal{X}_f and \mathcal{X}_c . The results demonstrate that when \mathcal{L}_{feat} is applied to lower *conv* layers, *DGCAN* can gener-

ate more distinct contour of the object from CAD-synthetic data and when \mathcal{L}_{coral} is applied to higher *conv* layers, *DGCAN* can generate more structured texture. Empirical evidence [9] shows that this effect mainly stems from two factors. ason is the increasing receptive field size, given the receptive field sizes of VGG-16’s *conv1_2*, *conv2_2*, *conv3_2*, *conv4_2*, *conv5_2* are 5, 14, 32, 76 and 164, respectively. The second factor is the increasing feature complexity along the network hierarchy.

To find the optimal trade off ratio λ , we synthesize images with λ ranging from 10^{-5} to 10^9 . The right plot in Figure 3 reveals when λ (\mathcal{L}_{coral} to \mathcal{L}_{feat} ratio) is small, the object contour will dominate the background texture cues. On the contrary, when λ is increased to 10^5 , the contour of the object gradually fades away and more structured textures from the real image emerge.

We randomly select some rendered results from three categories (“airplane”, “potted plant”, “sofa”) and show them in the left plot of Figure 4. The images are generated with the configuration $\mathcal{X}^f = conv3_2$, $\mathcal{X}^c = conv[1-5]_1$ ($\omega_c^{1\sim5} = 0.2$) and $\lambda = 10^3$. The results demonstrate that DGCAN-synthetic imagery preserves clear object contours from the CAD-synthetic images and synthesized textures from realistic domain.

We further leverage the DCNN visualization tool provided by [25] to show that DGCAN-synthetic images share more similarities with real images. [25] provides an effective tool to reconstruct an image from its representation. We compare the reconstruction results of bird images from three domains, i.e. DGCAN-synthetic, CAD-synthetic and real domains. In the right subplot of Figure 4, the odd rows show the original bird images, and their corresponding reconstructions are located in the even rows. From the plots,

²<http://shapenet.cs.stanford.edu/>

Method	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	overall
CAD-AlexNet [30]	25.7	19.8	11.8	31.3	11.4	72.0	26.0	8.4	12.2	27.7	1.0	3.2	3.3	69.1	11.6	40.0	0.0	15.4	43.7	40.4	18.48
CORAL-Alex [37]	33.1	15.9	17.7	27.7	3.5	79.9	26.5	17.6	15.0	22.8	5.7	10.0	11.1	62.1	12.3	29.1	0.0	9.1	25.5	26.6	18.18
DCORAL-Alex [38]	28.0	28.5	9.5	26.5	25.4	65.7	41.3	21.6	22.7	52.0	1.0	5.8	13.4	71.8	22.0	33.3	3.9	9.6	34.1	60.7	24.48
SA-fc7 [7]	59.8	46.5	34.2	29.3	6.2	66.9	28.4	31.4	15.6	23.1	9.7	11.7	19.7	63.7	10.3	29.9	9.0	24.5	16.2	46.5	21.10
GFK-fc7 [11]	43.4	31.4	18.6	41.7	4.6	65.7	24.5	10.3	16.8	14.9	15.1	4.9	18.0	50.7	5.2	23.1	12.9	17.2	14.9	29.4	16.14
DAN [20]	35.7	55.8	23.6	19.8	17.7	73.6	43.1	10.8	31.5	58.7	2.0	4.9	7.1	63.4	12.4	44.9	0.3	15.7	13.6	28.0	23.97
RTN [22]	20.9	3.6	2.6	12.0	1.5	69.7	37.2	33.8	21.8	53.2	0.0	0.2	0.3	61.0	11.6	21.5	0.0	3.5	7.6	44.6	17.76
StyleTransfer[9]	59.5	61.2	16.7	31.6	47.6	39.8	24.1	25.7	45.2	35.0	5.0	10.0	18.0	35.0	14.5	39.7	2.6	4.0	9.9	47.9	24.78
DGCAN-AlexNet	62.4	60.7	20.3	13.2	16.3	65.4	7.2	33.2	43.0	31.6	4.0	3.4	15.7	48.0	26.8	51.5	0.6	14.1	10.3	69.5	27.46

Table 1. **Results on PASCAL 2007.** We show per-category accuracy of different methods based on “AlexNet” [17] features. The results clearly demonstrate the superiority of our model over CAD-synthetic method [30], *style transfer* method [9] and several state-of-the-art domain adaptation models [37, 38, 7, 11, 20, 22]

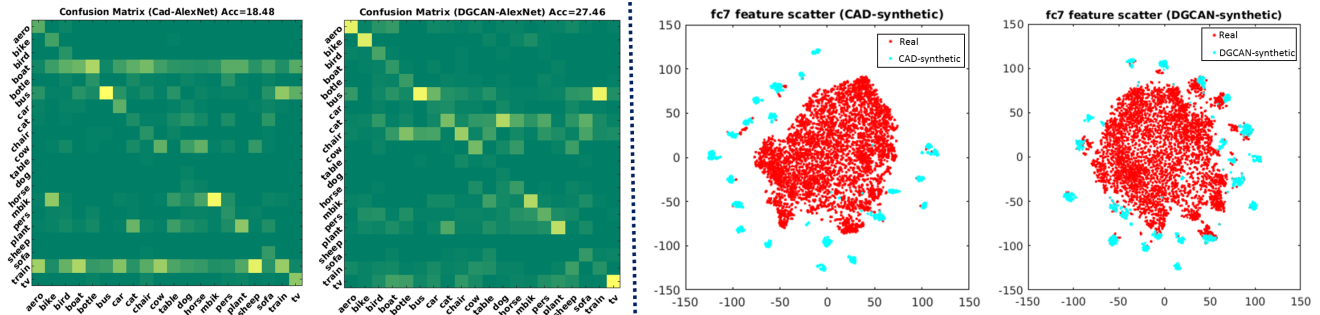


Figure 5. **Confusion Matrices and t-SNE plots of $fc7$ feature.** (1). The confusion matrices on the left show models trained on DGCAN-synthetic dataset (right subplot) pose a different error mode from those trained on CAD-synthetic dataset (left subplot). (2) The t-SNE plots on the right shows the embedded $fc7$ features of realistic and synthetic images are better aligned after applying our model to CAD-synthetic images. (Best viewed in color!)

	VGG	ResNet	AlexNet
CAD [30]	10.30	13.13	18.48
CORAL [37]	11.67	12.24	18.18
DCORAL [38]	17.76	-	24.48
SA [7]	20.38	19.33	21.10
GFK [11]	17.05	18.43	16.14
DGCAN	22.92	20.59	27.46

Table 2. **Results on PASCAL 2007.** We train three off-the-shelf classifiers on the adapted dataset and test on PASCAL 2007 benchmark. The results demonstrate the our model works better than other domain adaptation methods [38, 37, 7, 11] with all the three classifiers.

we can observe recognizable bird shapes from the reconstructed images of DGCAN-synthetic images. However, the birds in the reconstructed images of CAD-synthetic domain are lost in noisy color patches. These visualization results demonstrate that the DCNN can better recover category information from DGCAN-synthetic images than from CAD-synthetic images.

4.2. Domain Adaptation Experiments

In this section, we evaluate our approach on CAD-to-real domain adaptation tasks, using object classification as the application. The goal is to generate adapted CAD

images using our approach, then train deep object classifiers on the data, and test on real-image benchmarks. We compare the effectiveness of our model to previous methods [17, 38, 37, 11, 20, 22, 9, 7] on two benchmarks: PASCAL VOC 2007 [6] and the Office [32] dataset.

4.2.1 Experiments on PASCAL VOC 2007

Train/Test Set Acquisition As a training set, we generate 1080 images with DGCAN from the W-UG subset of CAD-synthetic dataset [30]. These images are equally distributed into 20 PASCAL categories. For evaluation, we crop 14976 patches from 4952 images in the test subset of PASCAL VOC 2007 dataset [6]. The patches are cropped using annotated object bounding boxes and each patch contains only one object.

Experimental Setup We evaluate the effectiveness of our approach by training three off-the-shelf DCNN classifiers, *i.e.* “AlexNet” [17], “VGG-16” [34] and “ResNet-50” (*residual net* with 50 layers) [14]. In the training process, the networks are initialized with the parameters pre-trained on ImageNet [4]. We replace the last output layer with a 20-way classifier and randomly initialize it with $\mathcal{N}(0, 0.01)$. We use mini-batch stochastic gradient descent (SGD) with a momentum of 0.9 to finetune all the layers. The base learning rate is 10^{-3} and the weight decay is 5×10^{-4} . Specifically, we set dropout ratios for $fc6$ and $fc7$ of “AlexNet” to

0.5. We report the results when the training iteration reaches 40k.

Baselines We compare our approach with the CAD-synthetic method [30], style transfer [9] and domain adaptation algorithms [37, 38, 7, 11, 20, 22].

To compare with state-of-the-art domain adaptation methods, we use the following baselines. **CORAL** [37] aligns the feature distribution of source domain ($P^s(\mathbf{x}^s, y^s)$) to target domain ($P^t(\mathbf{x}^t, y^t)$). **DCORAL** (Deep CORAL) [38] incorporates CORAL as a loss layer in the DCNN. **SA** (Subspace Alignment) [7] proposes a mapping function to align the subspace of the source domain with the target domain. The subspace is described by the eigenvectors of features [7]. **GFK** (Geodesic Flow Kernel) [11] models domain discrepancy by integrating numerous subspaces which characterize changes in geometric and statistical properties. Based on these subspaces, a geodesic curve is constructed, a geodesic flow kernel is computed and a kernel-based classifier is trained. **DAN** (Deep Adaptation Network) [20] and **RTN** (Residual Transfer Network) [22] train deep models with the Maximum Mean Discrepancy [33] loss to align the feature distribution of two domains.

For equal comparison, we take the same 1080 **W-UG** images which we utilized to generate our DGCAN-synthetic dataset as the source domain for domain adaptation algorithms. For style-transfer method [9], we use the same CAD-synthetic (*content*) images and real images (*style*) to generate new dataset. For **SA** [7] and **GFK** [11], we first extract deep features and then apply their model to get the baseline results. For all the baselines, we use the code and experimental settings provided by the authors to run all the experiments.

Results and Analysis The per-category accuracies of AlexNet classifier are presented in Table 1, demonstrating that our approach outperforms competing methods. After applying our approach to CAD-synthetic data, the overall accuracy rises from 18.48% to 27.46%. Additionally, Table 1 shows that our approach gains a clear advantage over the state-of-the-art domain adaptation algorithms [37, 38, 7, 11, 20, 22] and the style-transfer baseline [9]. The latter result reveals that aligning the covariance matrix works better than aligning the Gram matrix in the synthetic-to-real domain adaptation scenario. In Table 2, we further show VGG and ResNet classifiers trained on the adapted dataset outperform the CAD-synthetic method [30] and domain adaptation methods [37, 38, 7, 11]. With our model, the accuracies of VGG and ResNet classifiers rise from 10.3% to 22.92% and from 13.13% to 20.59%, respectively. We notice that AlexNet achieves the best overall performance. Given that VGG and ResNet have more parameters than AlexNet, we assume that they are overfitting to the generated synthetic dataset, which causes poor gen-

eralization to real-image domain.

We visualize how the inter-class confusion mode and the feature embeddings have changed after applying our model, as shown in Figure 5. The confusion matrices on the left show that AlexNet [17] trained on the CAD-synthetic dataset (\mathcal{I}_s) tends to mistake other categories for “boat” and “train”. This phenomenon disappears after applying our model to the CAD-synthetic dataset, as illustrated by the second subplot on the left of Figure 5. This effect is partially explained by the texture synthesizing ability of DGCAN, which provides additional discriminative cues to the CAD-synthetic images. At the feature level, we visualize layer *fc7*’s feature embeddings by t-SNE [23] before and after applying our model, as illustrated in the right two subplots of Figure 5. The t-SNE [23] visualization results clearly show that the features of realistic and synthetic images are better aligned after applying our model to CAD-synthetic images.

4.2.2 Experiments on the Office Dataset

We also evaluate our method on the Office benchmark [32], which was introduced specifically for studying the effect of domain shift on object classification. We evaluate the domain generalization ability of our approach by adapting the CAD domain to the real-image *Amazon* domain (images downloaded from amazon.com) in the Office dataset.

Train/Test Set Acquisition We apply our model to the 775 CAD-synthetic images provided by [30] to generate the training dataset. These CAD-synthetic images are rendered to train object detectors for Office dataset. To collect \mathcal{I}_t (natural images), for each category, we downloaded 4~5 images from Google by searching the category’s name. The test set comes from the Office Dataset [32], which has the same 31 categories (*e.g.* backpack, cups, *etc.*) in three domains, *i.e.* *Amazon*, *Webcam* (collected by webcam) and *DSLR* (collected by DSLR camera). Specifically, we use the *Amazon* set (2817 images) as the test set in our experiments as it is the most challenging setting, and because this domain significantly differs from PASCAL (see Table 4 for examples).

Baselines We compare our approach to two sets of baselines, with one set trained on another real image domain in Office (*Webcam* domain, 795 images) and the other trained on the CAD-synthetic domain (775 images). In both sets, we compare to the basic AlexNet [17] model (no adaptation) and domain adaptation algorithms [7, 11, 37, 20].

Results The results demonstrate that our approach performs strongly on this benchmark, as can be seen in Table 3. The overall classification accuracy of our model is 49.91% versus 44.69% for a classifier trained on the CAD-synthetic domain directly. The table also shows that *DGCAN* beats other baselines [7, 11, 37] and classifiers trained on real images (*Webcam* domain), and is slightly better than the domain

Method	bp	bk	bh	bc	bt	ca	dc	dl	dp	fc	hp	kb	lc	lt	mp	mt	ms	mg	pn	pe	ph	pr	pj	pn	rb	rl	sc	sp	st	td	tc	all
AlexNet-web [17]	76	96	90	48	33	86	82	59	3	51	60	70	76	14	14	60	86	32	32	55	42	52	30	12	24	43	41	16	31	39	14	47.30
SA-web [7]	77	96	87	27	34	77	84	35	5	40	63	71	73	8	29	47	77	33	32	42	35	57	67	14	26	50	37	13	13	34	17	47.18
GFK-web [11]	10	94	85	17	11	73	76	26	15	10	65	72	79	9	5	21	44	27	27	46	23	32	21	13	24	48	37	31	10	24	25	35.25
CORAL-web [37]	82	95	93	38	39	78	88	45	12	35	74	79	73	10	36	51	77	33	44	44	37	60	62	17	27	45	42	17	24	45	19	48.43
DAN-web [20]	87	96	82	23	58	87	90	35	2	36	78	85	74	9	57	89	83	68	43	48	35	52	44	17	37	43	42	28	19	29	16	49.63
AlexNet-CAD [17]	61	94	15	50	47	78	91	49	16	7	57	88	72	26	70	71	49	73	32	20	53	93	55	12	3	16	33	6	2	15	6	44.69
SA-CAD [7]	78	94	83	52	42	84	85	38	5	46	65	78	62	10	41	49	72	30	32	36	40	60	50	16	25	34	40	10	6	27	21	47.32
GFK-CAD [11]	60	96	78	17	19	69	84	22	18	16	49	56	68	12	1	23	35	22	20	45	30	12	16	14	24	35	42	8	6	21	23	31.38
CORAL-CAD [37]	60	98	92	63	39	79	90	42	15	36	80	80	74	12	37	55	63	33	32	47	46	58	60	19	30	39	44	20	16	38	19	47.93
DAN-CAD [20]	90	96	62	62	35	76	86	57	26	24	66	82	72	31	70	72	66	67	46	30	26	52	63	11	13	16	32	38	18	16	20	49.27
DGCAN	91	98	26	72	39	89	91	47	20	56	66	85	67	23	59	67	54	72	30	18	52	76	60	2	2	17	27	42	23	22	14	49.91
DGCAN+DAN	90	96	35	67	50	80	85	59	27	48	77	58	74	32	74	73	64	80	53	22	27	49	67	4	21	39	33	37	29	28	25	51.93

Table 3. **Results on Office Dataset.** We apply *DGCAN* to 775 CAD-synthetic images and train AlexNet classifiers [17]. The test images come from *Amazon* domain of Office dataset [32]. The results clearly shows our approach outperforms the competing baselines [37, 38, 7, 11, 20, 22]. The suffix “-web” and “-CAD” represent the methods are trained on *Webcam* domain [32] and CAD-synthetic domain, respectively.



















									
GT	back pack	bike helmet	bookcase	bookcase	bookcase	calculator	calculator	computer	file cabinet
Cad	printer	printer	bottle	computer	file cabinet	keyboard	phone	printer	bookcase
Ours	back pack	bike helmet	bookcase	bookcase	bookcase	calculator	calculator	computer	file cabinet
									
GT	keyboard	laptop	laptop	letter tray	mb phone	monitor	trash can	letter tray	desk lamp
Cad	computer	keyboard	notebook	monitor	calculator	monitor	trash can	printer	bike
Ours	keyboard	laptop	laptop	letter tray	mb phone	keyboard	ring binder	punchers	bike

Table 4. Instances from the Amazon domain of the Office dataset and the corresponding labels predicted by the baseline trained on CAD-synthetic images (CAD), and our model. We show examples where our model improves on the baseline, as well as typical failure cases.

alignment network (DAN) of [20]. We note here that this and other unsupervised domain adaptation baselines make use of the test data to train the alignment models (transductive training). On the other hand, our method did not use the test images for training, but performs well nonetheless.

We further show that our model is complementary with transductive domain adaptation algorithms. We set the newly generated dataset as the new source domain and adapt it to the real Amazon domain with *DAN* [20]. As showed in Table 3, this boosts the performance from 49.63% (49.27%) for *DAN* trained on real (CAD-synthetic) domain to 51.93%.

Table 4 shows some results for which the classifier trained on CAD-synthetic images fails to predict the correct labels while our model predicts the right ones, as well as some representative mistakes. The results show the potential to generate better training data for a large variety of object categories. An interesting example is the “desk lamp” with a toy bike in the middle, causing both models to mistake it for a bike.

5. Conclusion

Generating large-scale training examples from 3D CAD models is a promising alternative to expensive data annotation. However, the domain discrepancy between CAD-synthetic images and real images severely undermines the performance of deep learning models on real world applications.

In this work, we have proposed and implemented a Deep Generative Correlation Alignment Network to adapt the CAD-synthetic domain to realistic domains by generating images with more natural feature statistics. We demonstrated that leveraging ℓ^2 loss to preserve the content of the CAD models in feature space and applying the second-order *CORAL* loss to diminish the domain discrepancy are effective in synthesizing adapted training images. We empirically and experimentally show that *DGCAN*-synthetic images are more suitable for training deep CNNs than CAD-synthetic ones. An extensive evaluation on standard benchmarks demonstrates the feasibility and effectiveness of the proposed approach against previous methods. We believe our model can be generalized to other generic tasks such as pose estimation, saliency detection and robotic grasping.

References

- [1] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013. 2
- [2] K. M. Borgwardt, A. Gretton, M. J. Rasch, H.-P. Kriegel, B. Schölkopf, and A. J. Smola. Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006. 2
- [3] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. *arXiv preprint arXiv:1612.05424*, 2016. 3
- [4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009. 6
- [5] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 1
- [6] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>. 4, 6
- [7] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars. Unsupervised visual domain adaptation using subspace alignment. In *Proc. ICCV*, 2013. 1, 2, 6, 7, 8
- [8] L. A. Gatys, M. Bethge, A. Hertzmann, and E. Shechtman. Preserving color in neural artistic style transfer. *arXiv preprint arXiv:1606.05897*, 2016. 2
- [9] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015. 2, 3, 5, 6, 7
- [10] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. *arXiv preprint arXiv:1311.2524*, 2013. 1
- [11] B. Gong, Y. Shi, F. Sha, and K. Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2066–2073. IEEE, 2012. 1, 2, 6, 7, 8
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014. 2
- [13] R. Gopalan, R. Li, and R. Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *Proc. ICCV*, 2011. 2
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015. 1, 6
- [15] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *CoRR*, abs/1611.07004, 2016. 2
- [16] Y. Jia. Caffe: An open source convolutional architecture for fast feature embedding. <http://caffe.berkeleyvision.org/>, 2013. 4
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 1, 6, 7, 8
- [18] J. Liebelt and C. Schmid. Multi-view object class detection with a 3d geometric model. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 1688–1695. IEEE, 2010. 1, 2
- [19] M. Liu and O. Tuzel. Coupled generative adversarial networks. *CoRR*, abs/1606.07536, 2016. 2, 3
- [20] M. Long and J. Wang. Learning transferable features with deep adaptation networks. *CoRR*, abs/1502.02791, 1:2, 2015. 1, 2, 6, 7, 8
- [21] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu. Transfer joint matching for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1410–1417, 2014. 2
- [22] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Unsupervised domain adaptation with residual transfer networks. In *Advances in Neural Information Processing Systems*, pages 136–144, 2016. 1, 2, 6, 7, 8
- [23] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008. 7
- [24] A. Mahendran and A. Vedaldi. Understanding Deep Image Representations by Inverting Them. *ArXiv e-prints*, Nov. 2014. 2
- [25] A. Mahendran and A. Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *International Journal of Computer Vision*, pages 1–23, 2016. 5
- [26] F. Massa, B. C. Russell, and M. Aubry. Deep exemplar 2d-3d detection by adapting from real to rendered views. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6024–6033, 2016. 1, 2
- [27] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2
- [28] R. Nevatia and T. O. Binford. Description and recognition of curved objects. *Artificial Intelligence*, 8(1):77–98, 1977. 2
- [29] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. *arXiv preprint arXiv:1605.09304*, 2016. 2
- [30] X. Peng, B. Sun, K. Ali, and K. Saenko. Learning deep object detectors from 3d models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1278–1286, 2015. 1, 2, 3, 4, 5, 6, 7
- [31] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. In *European Conference on Computer Vision*, pages 102–118. Springer, 2016. 2
- [32] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *Computer Vision—ECCV 2010*, pages 213–226. Springer, 2010. 4, 6, 7, 8

- [33] D. Sejdinovic, B. Sriperumbudur, A. Gretton, K. Fukumizu, et al. Equivalence of distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics*, 41(5):2263–2291, 2013. [7](#)
- [34] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. [1](#), [3](#), [6](#)
- [35] M. Stark, M. Goesele, and B. Schiele. Back to the future: Learning shape models from 3d cad data. In *BMVC*, volume 2, page 5, 2010. [1](#), [2](#)
- [36] H. Su, C. R. Qi, Y. Li, and L. J. Guibas. Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2686–2694, 2015. [1](#), [2](#)
- [37] B. Sun, J. Feng, and K. Saenko. Return of frustratingly easy domain adaptation. *arXiv preprint arXiv:1511.05547*, 2015. [1](#), [2](#), [4](#), [6](#), [7](#), [8](#)
- [38] B. Sun and K. Saenko. Deep CORAL: correlation alignment for deep domain adaptation. *CoRR*, abs/1607.01719, 2016. [1](#), [2](#), [4](#), [6](#), [7](#), [8](#)
- [39] M. Sun, H. Su, S. Savarese, and L. Fei-Fei. A multi-view probabilistic model for 3d object classes. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 1247–1254. IEEE, 2009. [1](#), [2](#)
- [40] E. Tzeng, C. Devin, J. Hoffman, C. Finn, X. Peng, S. Levine, K. Saenko, and T. Darrell. Towards adapting deep visuo-motor representations from simulated to real environments. *arXiv preprint arXiv:1511.07111*, 2015. [2](#)
- [41] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. [2](#)
- [42] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014*, pages 818–833. Springer, 2014. [2](#)