

Representation Learning Through Latent Canonicalizations

Or Litany
NVIDIA*

Ari Morcos
Facebook AI Research

Srinath Sridhar
Brown University

Leonidas Guibas
Stanford University*

Judy Hoffman
Georgia Tech*

Abstract

*We seek to learn a representation on a large annotated data source that generalizes to a target domain using limited new supervision. Many prior approaches to this problem have focused on learning “disentangled” representations so that as individual factors vary in a new domain, only a portion of the representation need be updated. In this work, we seek the generalization power of disentangled representations, but relax the requirement of explicit latent disentanglement and instead encourage linearity of individual factors of variation by requiring them to be manipulable by learned linear transformations. We dub these transformations **latent canonicalizers**, as they aim to modify the value of a factor to a pre-determined (but arbitrary) canonical value (e.g., recoloring the image foreground to black). By creating simple simulators with pre-specified factors of variation to roughly approximate datasets such as SVHN and ImageNet, we demonstrate experimentally that our method helps reduce the number of observations needed to generalize to a similar target domain when compared to a number of supervised baselines.*

1. Introduction

Most state-of-the-art visual recognition models rely on supervised learning using a large set of manually annotated data (21; 11; 10; 38). As recognition task complexity increases, so does the number of potential real world variations in visual appearance and hence the size of the example set needed for sufficient test time generalization. Unfortunately, large labeled data sets are laborious to acquire (5; 51), and may even be infeasible for applications with evolving data distributions.

Often a large portion of the variance within a collection of data is due to task-agnostic factors of variation. For example, the appearance of a street scene will change substantially

based on the time of day, weather pattern, and number of traffic lanes, regardless of whether cars or pedestrians are present. Ideally, the ability to recognize cars and pedestrians would not require labeled examples of street scenes for all combinations of times of day, weather conditions, and geographic locations. Rather it should be sufficient to observe examples from each factor independently and generalize to unseen combinations. However, often the in-domain labeled data available may not even linearly cover all factors of variation. This calls for methods that encourage such sample efficiency by focusing on the individual complexities of the factors of variation, as opposed to their product.

In this work, we propose learning a factored representation by leveraging a large collection of source domain data with meta-labels specifying the factors of variation within an image. Such a collection may be available from meta-data, attribute labels, or from hyper-parameters used for generation of simulated imagery - which, we show experimentally, need not be a realistic rendering of the target domain data and can be easily approximated for common datasets such as SVHN and ImageNet. Prior approaches to learn from a source domain and ignore factors of variation consider learning domain invariant representations (8; 45). While, prior approaches to learning representations which isolate factors of variation in the data have typically regularized the representation itself, with the aim of learning “disentangled” representations (22; 3; 12; 13; 1; 17).

In this work, we propose using the source data with known factors of variation to regularize the *way the representation can be manipulated* rather than the representational structure itself. Here, we take such an approach by introducing **latent canonicalization**, in which we constrain the representation such that individual factors can be clamped to an arbitrary, but fixed value (“canonicalized”) through a simple linear transformation of the representation. These canonicalizers are learned such that they can be applied independently or composed together to canonicalize multiple factors. Latent canonicalizers are optimized by a pixel loss

*Majority of work done while at Facebook AI Research

over pairs of ground-truth canonicalized examples and reconstructions of images with various factors of variation whose representations have been passed through the relevant latent canonicalizers. By requiring the ability for manipulation of the latent space according to factors of variation, latent canonicalization encourages the linearization of such factors.

We evaluate our approach on its ability to learn general representations after observing only a subset of all potential combinations of factors of variation. We first consider the simple dSprites dataset, introduced to study disentangled representations (30) and show qualitatively that we can effectively canonicalize individual factors of variation. We next consider the more realistic, though still tractable, task of digit recognition on street view house numbers (SVHN) (34) with few labeled examples. Using a simulator we designed to roughly approximate SVHN, we train our representation with latent canonicalization along multiple axes of variation such as font type, background color, etc. and then use the factored representation to enable more efficient few-shot training of real SVHN data. Our method substantially increased performance over standard supervised learning and fine-tuning for equivalent amounts of data, achieving digit recognition performance that was only attainable with $5\times$ as much SVHN labeled training data for the best-performing baseline method. Finally, to demonstrate that our approach scales to naturalistic images, we evaluate our method on a subset of ImageNet using a simulator constructed from ShapeNet (2), again outperforming the best baselines. Our experiments offer promising evidence that encoding structure into the latent representation guided by known factors of variation within data can enable more data efficient learning solutions.

2. Related Work

2.1. Sim2Real

The setting of near-unlimited simulated data with ground truth labels and scarce real data occurs often in computer vision and robotics. However, the *domain gap* between simulated and real data reduces generalization capacity. Many approaches have been proposed to overcome this difficulty which are broadly referred to as *sim2real* approaches. A simple approach to closing the *sim2real* gap is to train networks with combinations of real and synthetic data (47).

Transfer Learning and Few-shot Learning: Simulated data is most useful when there is a shortage of labeled real data (39). In this situation, one may make use of few-shot learning techniques which seek to prevent over-fitting by using a metric loss between data triplets (20), comparing similarity between individual examples (50) or between a prototypical class example and each instance (42). Other techniques use meta-learning approaches (7). (33) modeled object attributes as learned operators on object vectors,

though that work primarily focused on compositional generalization.

Domain adaptation: With access to a large set of unlabeled real examples, domain adaptation techniques can be used to close the *sim2real* gap. One class of techniques focuses on matching domains at the feature level (8; 29; 45), aiming to learn domain-invariant features than can make models more transferrable (53; 23). In fact, image-to-image translation focuses on the appearance gap by bridging the *appearance gap* in the image domain instead of feature space (41; 27). Domain adaptation can also be used to learn *content gap* in addition to appearance gap (16). Additional structural constraints, such as cycle consistency, can further refine this image domain translation (52; 14; 32). Finally, image stylization methods can also be adapted for *sim2real* adaptation (24).

Domain randomization: (DR) exploits control of the synthetic data generation pipeline to randomize sources of variation (44; 36). Random variations will likely be ignored by networks and thus result in invariance to those variations. A particularly interesting instantiation of DR was suggested in (43) for pose estimation. Pose is an example of a factor of variation which could be ambiguous due to occlusions and symmetries. Instead of explicitly regressing for the pose angle, the authors propose an implicit latent representation. This is achieved by an augmented-autoencoder, a form of denoising-autoencoder that addresses all nuisance factors of variation as noise. This idea can be seen as a version of our method in which all factors are canonicalized at once rather than individually. Another interesting example is the quotient space approach of (31) which removes pose information for a 3D shape representation by max-pooling encoder features over a sampling of object rotations. It, however, does not consider how to perform canonicalization as a linear transformation in latent space, nor how to compose different canonicalizers.

2.2. Disentangling

A number of studies have sought to learn low-dimensional representations of the world, with many aiming to learn “disentangled” representations. In disentangled representations, single latent units independently represent semantically-meaningful factors of variation in the world and can lead to better performance on some tasks (46). This problem has been most commonly studied in an unsupervised setting, often by regularizing latent representations to stay close to an i.i.d Gaussian prior (12; 13; 1; 17). An extension to long-tail distributions was shown in (18). Other popular unsupervised approaches include maximizing the mutual information between the latents and the observations (3) and adversarial approaches (6; 37) and When supervision on the sources of variation is available, it is possible to use this in a weak way (22). Disentanglement to primitives functions

was also studied in the context of compositional generalization (26; 25).

Many of these works have explicitly endeavored to learn semantically meaningful representations which are both linearly independent and axis-aligned, such that individual latents correspond to individual factors of variation in the world. However, recent work has questioned the importance of axis-aligned representations, showing that many of these methods rely on implicit supervision and finding little correspondence between this strict definition of disentanglement and learning of downstream tasks (40; 28). Further, while axis-alignment is useful for human interpretability, for the purposes of decodability, any arbitrary rotation of these latents would be equally decodable so long as factors are linearly independent (28). In this work, we use explicit supervision in a simulated setting to encourage linear, but not necessarily axis-aligned representations.

3. Approach

Our goal is to learn representations which capture the generative structure of a dataset by independently representing the underlying factors of variation present in the data. While many previous works have approached this problem by regularizing the representation itself (22; 3; 12; 13; 1; 17), here we take a different approach: rather than directly encourage the representation to be disentangled, we instead encourage the representation to be structured such that individual factors of variation can be manipulated by a simple linear transformation. In other words, we constrain the *way that the representation can be manipulated* rather than the *structure of the representation itself*.

3.1. Latent canonicalization

In our approach, we augment a standard convolutional denoising autoencoder (AE) with **latent canonicalizers**. A standard AE learns an encoder, Enc , which takes as input a given image, x , and produces a corresponding latent vector, z . At the same time, the latent vector is used as input to a decoder, Dec , which produces an output image, \hat{x} . Both the encoder and decoder are learned according to the following objective, \mathcal{L}_{ae} , which minimizes the difference between the original input image and the reconstructed output image:

$$z = \text{Enc}(x; \theta_e) \quad ; \quad \hat{x} = \text{Dec}(z; \theta_d); \quad (1)$$

$$\mathcal{L}_{\text{ae}}(\theta_e, \theta_d; x) = \min_{\theta_e, \theta_d} \|x - \hat{x}\|_2^2. \quad (2)$$

To encourage noise-robustness, we augment the potential input images following previous work on denoising autoencoders (48; 49), noising each raw input image, x , by adding Gaussian noise, blur, and random crops and rotations, leading to our noised input image, \tilde{x} .

In this work, we additionally constrain the structure of the learned latent space using a set of latent canonicalization

losses. We define a latent canonicalizer as a learned linear transformation, \mathcal{C} , which operates on the latent representation, z , in order to transform a given factor of variation (e.g., color or scale) from its original value to an arbitrary, but fixed, canonical value. So that individual factors can be manipulated separately, we learn unique canonicalization matrices, \mathcal{C}_j , for each factor of variation, $j \in [1, K]$, present in the dataset. In order to constrain the latent representation according to canonicalization along one factor, j , our method yields the following basic form:

$$z_{\text{canon}}^{(j)} = \text{Enc}(\tilde{x}; \theta_e) \cdot \mathcal{C}_j \quad ; \quad \hat{x}_{\text{canon}}^{(j)} = \text{Dec}(z_{\text{canon}}; \theta_d) \quad (3)$$

To supervise the learning of latent canonicalizers, we compare the images generated by canonicalized latents, \hat{x}_{canon} , to ground truth images with the appropriate factors of variation set to their canonical values, x_{canon} . Canonicalizers can also be composed together to canonicalize multiple factors (e.g., $z_{\text{canon}}^{(j,k)} = z \cdot \mathcal{C}_j \cdot \mathcal{C}_k$). During training, each image is passed through a random subset of both individual and pairs of canonicalizers. Given Q canonicalization paths for a given image, x , the corresponding canonicalization loss for that single image (red in Figure 1) is written as:

$$\mathcal{L}_{\text{canon}} = \frac{1}{Q} \sum_q \|\hat{x}_{\text{canon}}^{(q)} - x_{\text{canon}}^{(q)}\|_2^2 \quad (4)$$

Since many outputs are canonicalized, it is possible that the decoder will simply learn to only generate the canonical value of a given factor of variation. To prevent this form of input-independent memorization, we also include a “bypass” loss which is equivalent to the standard denoising auto-encoder formulation (green in Figure 1) defined in Equation (2), thus forcing information about each factor to be captured in the latent vector, z .

Finally, we ensure that our representation not only allows for linear manipulation along factors of variation, but does so while capturing the information necessary to train a classifier to solve our end task. To this end, we add a supervised cross-entropy loss, \mathcal{L}_{CE} , which optimizes our end task using available labeled data (cyan in Figure 1):

$$\mathcal{L}_{\text{CE}} = y \log \hat{y} \quad (5)$$

Combining equations 2, 4, and 5 with loss-scaling factors (α and β) gives us our final per-example loss formulation:

$$\mathcal{L} = y \log \hat{y} + \alpha \|\hat{x} - x\|_2^2 + \beta \frac{1}{Q} \sum_q \|\hat{x}_{\text{canon}}^{(q)} - x_{\text{canon}}^{(q)}\|_2^2 \quad (6)$$

In practice, two canonicalizers are chosen at random for each input batch, \mathcal{C}_h and \mathcal{C}_j , and the corresponding latent representation z is passed through $\{\mathcal{C}_h, \mathcal{C}_j, \mathcal{C}_h \mathcal{C}_j, \mathcal{C}_j \mathcal{C}_h\}$ generating four unique canonicalized latents: $\{z_{\text{canon}}^{(h)}, z_{\text{canon}}^{(j)}, z_{\text{canon}}^{(h,j)}, z_{\text{canon}}^{(j,h)}\}$. A diagram of

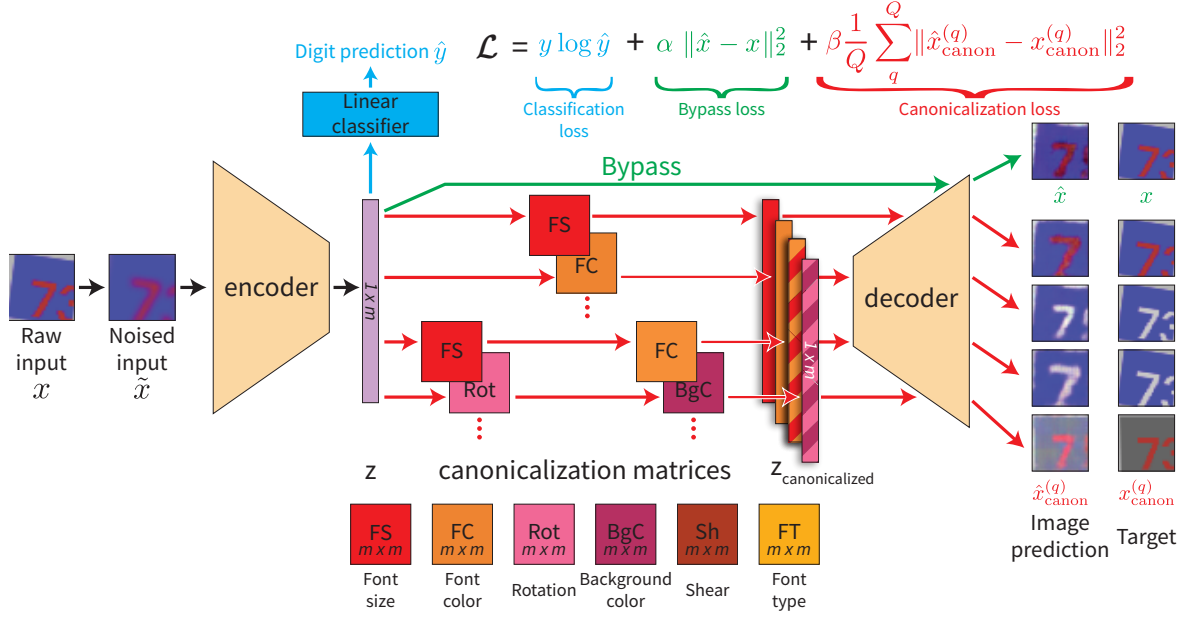


Figure 1: Schematic representation of latent canonicalization: Colored paths correspond to different components of the loss (cyan: classification, green: bypass, red: canonicalization). Four possible canonicalizations (two individual and two pairs) are shown along with example simulated SVHN images and reconstructions.

the method is shown in Figure 1, illustrating each canonicalization path, the bypass path, and the classification model. We have thus far focused on the single image loss for simplicity. The full model averages the per image loss over mini-batch before making gradient updates.

Latent canonicalizer constraints: Our approach relies upon constraining the way a representation can be manipulated. As a result, the specific choice of constraints should have a significant impact on the representations which are ultimately learned. Here, we limit ourselves to only two constraints: the transformations must be linear and canonicalizers must be composable, at least in pairs. If we were to allow non-linear canonicalizers, there would be little incentive for the encoder to learn an easily manipulable embedding. This would only be exacerbated as the non-linear canonicalizer is made more powerful by e.g., additional depth. By requiring canonicalizers to be composable, we encourage independence as each canonicalizer must be able to be applied without damaging any other. We explore some other potential constraints in Section 5.

3.2. sim2real evaluation

A main motivation of latent canonicalization is to leverage structure gleaned from a large source of data with rich annotations to better adapt to downstream tasks. A natural such setting for performance evaluation is sim-to-real. Specifically, we make use of the Street View House Numbers (SVHN) dataset and a subset of ImageNet (5) as our real

domains. To simulate SVHN, we built a SVHN simulator in which we have full control over many factors of variation such as font, colors and geometric transformations (a detailed description of the simulator is given in Section 4.2.1). To simulate ImageNet, we built a simulator which renders 3D models from ShapeNet (2) to generate ImageNet-like images (see Section 4.3.1 for details).

We first pre-train on the synthetic data with latent canonicalization. Following this step, we freeze the canonicalizers and investigate whether the learned representations can be leveraged as the input to a linear classifier for few-shot learning on real examples, labeled only with the class of interest (e.g., no meta-labels for additional factors of variation). During this stage, the encoder is also refined.

Majority vote: Because latent canonicalization manipulates the latent representation, we can use canonicalization as a form of “latent augmentation.” In this setting, we can aggregate the predictions of the digit classifier across many canonicalization paths, each of which confers a single “vote.” Critically, such an approach requires the ability to cleanly manipulate the learned representation, and is therefore only possible for our proposed method. For a more detailed exploration of the impact of majority vote, see Section 5.

3.3. Baselines

We compare our proposed latent canonicalization with several baselines. For fair comparison, we fixed as many hyperparameters as possible: we use the same backbone architecture in all our networks (details are in Section 3.1);

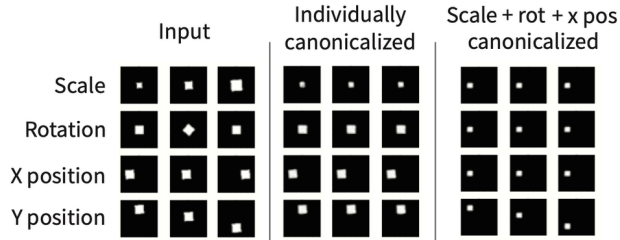


Figure 2: Canonicalization of dsprites images: Input dsprites images (left), reconstructions of inputs with one factor canonicalized (middle), and rotation, scale, and x-position canonicalized (right). Each row shows how images change as a single factor of variation is altered.

the same number of epochs for pre-training (learning from synthetic data); and a carefully chosen number of epochs at the refinement stage to fit well the method overfitting rate. For all latent canonicalization experiments we trained three models on the simulated data and then performed five refinements of each pre-trained model, for a total of 15 replicates. Results are reported as mean \pm standard deviation. For baselines, 15 independent replicates were trained.

Our simplest baseline, which is meant as a lower bound, is simply a classifier trained only on the low-shot real data. For pre-training methods, we compare to two categories of baselines: purely unsupervised methods and methods which have access to supervision on the simulated data. Among unsupervised methods, we compare against simple deterministic and variational autoencoders pre-trained on synthetic data, after which a linear classifier is trained on low-shot real data (Vanilla AE and VAE, respectively), the beta-VAE model with two values of beta (12), and a self-supervised rotation prediction task (9). We include these comparisons for context and completeness, but emphasize that these methods *do not* have equivalent supervision to our method. A better comparison is to models which have equivalent access to supervision on the synthetic data. We therefore compare to a classifier pre-trained on synthetic data and refined on low-shot real data and to a deterministic and variational autoencoder pre-trained with a digit classifier during both pre-training (on synthetic data) and training on few-shot real data. Not surprisingly, these are generally the strongest of our baselines. The loss weighting was determined individually for each model. Finally, to measure the importance of using constrained, linear canonicalizations, we also include a latent canonicalizer baseline that does not impose a linearity constraint on the latent space, i.e., we replace the learned linear transformation \mathcal{C} by a 2-layer MLP.

4. Experiments

4.1. Latent canonicalization of dSprites

Key to our method is the use of latent canonicalizers, which are learned linear transformations optimized to eliminate individual factors of variation. As a first test of the effectiveness of latent canonicalization, we evaluated our

framework using the toy data set, dSprites (30), which was designed for the exploration and evaluation of disentanglement methods. Specifically, dSprites is a dataset of images of white shapes on a black background generated by five independent factors of variation (shape, scale, rotation, and x- and y-positions). Training our model (Figure 1) on dSprites, we demonstrate the effect of applying different individual canonicalizers to various values of the input factors (Figure 2 left and middle). We therefore also applied a set of three canonicalizers (scale, rotation, and x-position) sequentially as shown in Figure 2, right. Encouragingly, we found that not only did individual canonicalizers effectively canonicalize their factor of interest, multiple canonicalizers can be applied in sequence. Furthermore, although models were trained with only pairs of canonicalizers, triplets of canonicalizers also performed well (Figure 2, right).

4.2. Latent canonicalization of SVHN

4.2.1 Simulating SVHN

To support our proposed training procedure, we require a comprehensive dataset with detailed meta-data about ground-truth factors of variation. While this is possible for a natural dataset, such data can also be generated for visually realistic, but fairly simplistic datasets such as SVHN. To this end, we built a procedural image generator that simulates the SVHN dataset by rendering images with digits on a constant-colored background (see Figure A1). Apart from digit class variation, we also simulate: font color, background color, font size, font type, rotation, shear, fill color for newly created pixels, scale, number of digit instances, translation, Gaussian noise, and blur. A detailed description is provided in Section B. Among these factors we chose the first six for supervision, noise and blur as a joint noise model, and the rest as additional factors to enrich the data variety without supervised canonicalization. Some of the resulting images can be seen in Figure 3 (see Appendix B for further simulator details). To enable reproducibility across comparisons, and to minimize unaccounted for variability in the data, we generated a fixed training set with 75,000 images along with targets for all possible canonicalization paths. We emphasize that this training set represents a small fraction ($\sim 0.2\%$) of the total number of possible combinations of factors. We used such a small fraction of the total space to demonstrate that latent canonicalization is feasible even if the factor space is only sparsely sampled. The simulator along with the generated train set will be made publicly available.

4.2.2 sim2real SVHN transfer using latent canonicalization

We want to learn representations which enable models to generalize to novel data with consistent underlying structure.

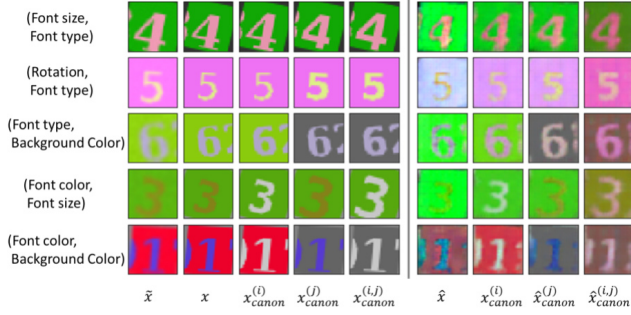


Figure 3: Example targets and reconstructions of canonicalized simulated SVHN images.

Model	10 shot	20 shot	50 shot	100 shot	1000 shot
Vanilla AE	15.39 ± 6.12	21.55 ± 9.16	41.98 ± 16.62	51.18 ± 19.80	58.51 ± 19.68
VAE	15.73 ± 2.48	22.02 ± 2.76	44.57 ± 4.34	68.73 ± 2.41	84.41 ± 0.52
beta-VAE (12) ($\beta = 5$)	12.09 ± 0.67	13.33 ± 0.23	17.20 ± 1.10	25.85 ± 3.02	56.23 ± 3.28
beta-VAE (12) ($\beta = 10$)	12.20 ± 0.95	13.53 ± 0.50	15.10 ± 1.86	19.47 ± 1.41	35.82 ± 4.90
RotNet (9)	49.25 ± 1.04	66.42 ± 2.16	79.38 ± 0.5	86.26 ± 1.11	90.58 ± 0.17
Classifier (real only)	20.86 ± 2.03	36.08 ± 2.49	72.75 ± 1.99	82.24 ± 0.78	92.84 ± 0.39
Classifier (synth only)	76.50 ± 2.06	80.07 ± 1.09	83.95 ± 0.79	85.65 ± 1.17	89.82 ± 0.60
AE + Classifier	79.61 ± 1.18	81.63 ± 1.00	84.66 ± 0.76	85.86 ± 0.82	88.80 ± 0.62
VAE + Classifier	78.46 ± 0.18	80.77 ± 0.71	84.49 ± 1.00	86.16 ± 0.30	90.05 ± 0.59
Our (nonlinear C)	79.13 ± 0.28	81.61 ± 0.67	85.18 ± 0.72	87.01 ± 0.42	89.98 ± 0.35
Ours (linear C)	82.55 ± 0.86	84.83 ± 0.76	87.82 ± 0.57	89.40 ± 0.48	91.21 ± 0.24
Ours + majority vote	83.41 ± 1.23	85.41 ± 0.88	88.17 ± 0.53	89.58 ± 0.57	91.34 ± 0.34

Table 1: SVHN sim2real transfer results Model performance on the SVHN test set using low-shot labeled real examples for method and baselines. Table entries represent mean ± std.

Moreover, the effectiveness of disentangling for acquisition of downstream tasks has recently been called into question (28). We therefore evaluate the quality of our learned representations by measuring their ability to adapt to real examples. Specifically, we consider a few-shot setting where models pre-trained on simulated data have access to a small refinement set of a few annotated examples per class. We ran this experiment with per-class set sizes of 10, 20, 50, 100 and 1000. To measure sim2real transfer, we train a fresh linear classifier on the representation learned by the encoder pre-canonicalization, z , while also allowing the encoder to be refined. We report accuracy on the unseen SVHN test set. As a measure of the pre-trained model, Figure 3 includes examples of reconstructions generated by canonicalized latents. When small train-sets are used, results may vary substantially depending on the selection. To account for this, we (a) use the same train set across methods, and (b) ran each experiment 15 times: 3 different networks were trained on simulated data with different random seeds and 5 replicate refinements were performed per pre-trained network.

Table 1 shows the sim2real SVHN results for our method with eight baselines discussed in Section 3.3. For all settings with fewer than 1000 examples per class, we found that our model outperformed the best competing baseline by ~ 3 -4%. Notably, the use of constrained, linear canonicalizers

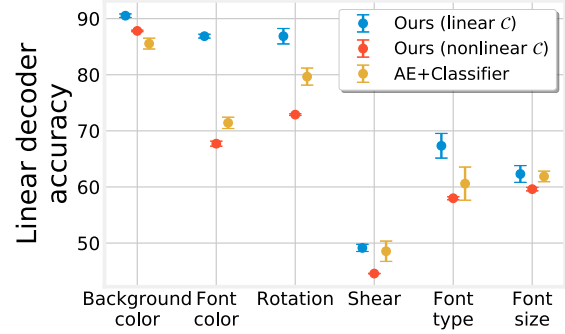


Figure 4: Linear decodability of factors of variation. Performance of a linear classifier trained on the frozen, pre-canonicalized representation, z , for each factor of variation. Chance is 1.6% for background and font color, 33.3% for rotation, 10% for shear, and 16.7% for font type and size. Error bars represent mean ± std across three pre-trained networks.

substantially improved performance over more expressive, nonlinear canonicalizers, highlighting the importance of constrained manipulations of latent space. We further improved performance by taking a “majority vote” approach, in which we pass the representation through multiple latent canonicalizers in parallel to generate multiple votes as discussed in Section 3.2. Consistently, we found that majority vote boosted performance, by up to $\sim 0.9\%$, with the largest gains coming for the lowest-shot settings. In Section 5, we discuss four additional directions for improvement which, unfortunately, either harmed or left unchanged sim2real performance relative to our best-performing models.

To contextualize the importance of this improvement, one can see that to match our reported performance on 20 shot with the best baseline of an AE + Classifier, a 5 times larger train set of 100 is required. This demonstrates the potential of our proposed method in better utilizing access to meta-labels for better adaptation to real data.

4.2.3 Analysis of representations

Linear decodability of factors of variation from representations: While latent canonicalization encourages representations to be linearly manipulable, it does not explicitly encourage linear decodability. However, since our canonicalizers are constrained to be linear, latent canonicalization may also encourage linear decodability. To test this, we trained linear classifiers on the pre-trained, frozen encoder for each factor of variation. We ran this experiment separately for each factor and compared linear decodability to our best baseline, AE+Cls (yellow), and to models trained with nonlinear canonicalizers (red; to measure the importance of constrained, linear canonicalizers). For background color, font color, and rotation angle, the canonicalized representation was noticeably more linear than the baseline (shown here by higher accuracy on a held-out test set), whereas font type showed a smaller improvement and font size and shear showed no improvement in linear decodability (Figure 4). One possible explanation for the discrepancy across

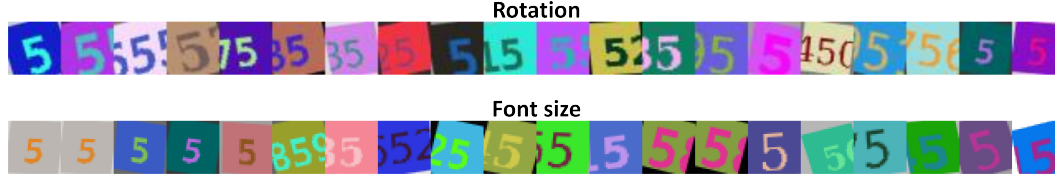


Figure 5: Linear properties of the representation. Each row shows the first principal component of $\mathbf{z}_{\text{canon}}^{(j)} - \mathbf{z}$ for a source of variation. A clear pattern is visible for rotation (left to right tilted), and font size (small to big). 20 normally distributed samples from a batch of 1000 are shown above. See supplementary for more visualizations.

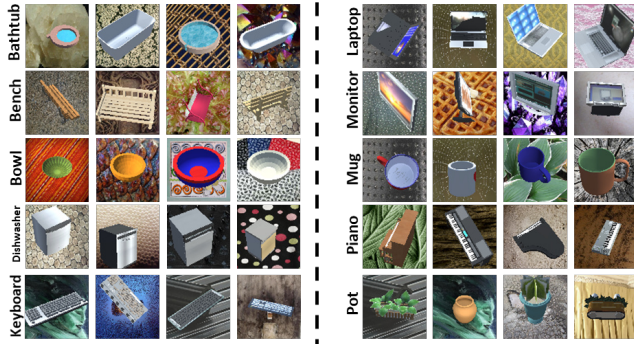


Figure 6: SynthImageNet dataset

factors is that font color, background color and, rotation are the most visually salient factors with the largest range of variability. Critically, linear canonicalization outperformed nonlinear canonicalization for all properties, demonstrating the importance of using constrained, linear canonicalizers.

Visualizing the impact of canonicalization: If these representations were indeed linear, we would expect them to be easily decomposable using principal component analysis (PCA), the components of which we can visualize. However, the latent codes from each of the canonicalizers *remove* the effect of a source of variation while keeping the others. We therefore compute the principal components (PCs) of $\mathbf{z}_{\text{canon}}^{(j)} - \mathbf{z}$, i.e., the difference between a canonicalized latent and the pre-canonicalized latent, such that PCs now represent the *removed* factor of variation. In Figure 5, we show sorted images along the first PC, showing a clear linear sorting of rotation, and font size. This visually demonstrates that our approach is able to extract latents that have strong linearity.

4.3. Latent canonicalization of ImageNet subset

4.3.1 Simulating ImageNet: *SynthImageNet*

To demonstrate the few-shot sim2real transfer capability of our method on a more naturalistic, complex dataset, we built a simulator to synthesize images similar to ImageNet (5). Our simulator uses 3D models from ShapeNet (2) to render plausible images of different shapes from various camera

Model	10 shot	20 shot	100 shot
Vanilla AE	19.27 ± 4.44	24.60 ± 2.25	34.33 ± 2.00
VAE	19.47 ± 0.64	23.87 ± 1.86	33.13 ± 1.55
β -VAE (12) ($\beta = 5$)	18.87 ± 0.7	22.87 ± 1.63	31.93 ± 1.81
β -VAE (12) ($\beta = 10$)	16.20 ± 1.40	19.00 ± 2.43	28.93 ± 2.00
RotNet (9)	22.07 ± 0.42	26.13 ± 1.42	37.13 ± 0.31
Classifier (real only)	23.13 ± 1.50	29.07 ± 3.01	38.27 ± 4.00
Classifier (synth only)	36.00 ± 1.91	38.87 ± 0.99	45.13 ± 0.76
AE + Classifier	33.93 ± 0.58	37.07 ± 3.23	43.07 ± 1.86
VAE + Classifier	34.73 ± 0.90	37.73 ± 0.83	44.27 ± 0.81
Ours (nonlinear \mathcal{C})	35.33 ± 0.50	37.47 ± 1.15	44.40 ± 1.22
Ours (linear \mathcal{C})	39.66 ± 1.40	40.84 ± 1.36	46.07 ± 2.12
Ours + majority vote	39.00 ± 0.72	40.47 ± 1.86	46.00 ± 0.80

Table 2: sim2real transfer on ImageNet subset Model performance on the 10 class ImageNet test set using low-shot labeled real examples for method and baselines. Table entries represent mean \pm std.

orientations and scales (Figure 6). To evaluate few-shot transfer from simulated to real ImageNet, we chose a subset of 10 classes which overlapped with ShapeNet categories (“ImageNet subset”). For each class, we rendered a total of 5000 frames, each containing a randomly chosen 3D model instance from the category. To increase variability, we also augment the background of each image with a randomly chosen texture from the Describable Textures dataset (4). We consider 4 factors of variation for this synthetic dataset, which we call *SynthImageNet*: camera orientation (latitude, longitude), object scale, and background texture.

4.3.2 sim2real ImageNet subset transfer using latent canonicalization

Table 2 shows sim2real results on the 10-class subset of ImageNet. Our method shows consistent improvement over baselines demonstrating that latent canonicalization can generalize to more naturalistic and complex settings.

5. Alternative design decisions

Latent canonicalization opens up many additional avenues for modification to potentially produce better representations and, consequently, better sim2real performance.

In the previous section, we showed how incorporating majority vote further increased performance. Here, we discuss several other modifications we explored, which resulted in no change or a decrease in sim2real performance.

Idempotency reconstruction loss: To encourage composability of canonicalizers, we trained them in identical pairs with alternating orders for consistency (e.g., $z \cdot C_1 C_2$ and $z \cdot C_2 C_1$). We also tried encouraging idempotency, such that the same canonicalizer can be repeatedly applied without changing the reconstruction (e.g., $z \cdot C_1 C_1$). We found that applying this loss actually *harmed* performance, reducing pre-majority vote classification accuracy by $\sim 0.5\%$ (Table 3, second row).

Classifier location during pre-training: In our model, the classifier is placed at the output of the encoder prior to canonicalization, z . However, one might imagine that latent canonicalization could serve as a form of data augmentation, such that placing the classifier after the canonicalization step, z_{canon} Canon, might increase performance. In contrast, we found that placing the classifier after the canonicalization step harmed performance, reducing pre-majority vote classification accuracy by $\sim 3\%$ (Table 3, third row). Interestingly, however, we found that the majority vote method, which can also be viewed as a form of data augmentation, did in fact increase performance.

Latent consistency and idempotency loss: The impact of latent canonicalization was supervised at the image level, by comparing the reconstruction to a target image. However, we could also use a self-supervised loss at the latent level, by enforcing consistency (i.e., $\min \|z \cdot C_1 C_2 - z \cdot C_2 C_1\|_2$) and idempotency (i.e., $\min \|z \cdot C_1 C_1 - z \cdot C_1\|_2$). To account for the scale of this latent loss, which was much larger than the other loss components, we used a very small scale factor for the latent loss to maximize performance ($1e-7$). Even with an appropriately scaled loss, however, the latent loss either had little impact or harmed sim2real performance (Table 3, fourth row).

Alternative majority votes: We found that a simple majority vote containing the pre-canonicalized and individually canonicalized representations (1 and 6 votes, respectively) increased performance by $\sim 0.5\%$. To further augment the vote set, we also tried adding votes via idempotency (6 additional votes) and pairs (30 additional votes). We found that neither of these additions further improved performance over the simplest majority vote approach (Table 4).

Model Variants	10 shot	20 shot	50 shot	100 shot	1K shot
No maj vote	82.55 \pm 0.86	84.83 \pm 0.76	87.82 \pm 0.57	89.40 \pm 0.48	91.21 \pm 0.24
+ idem	81.77 \pm 1.23	84.08 \pm 0.98	87.08 \pm 0.64	88.96 \pm 0.51	90.74 \pm 0.35
+ classifier after	79.69 \pm 1.22	81.14 \pm 1.00	84.21 \pm 0.74	86.42 \pm 0.54	88.62 \pm 0.24
+ latent loss $1e-7$	82.74 \pm 0.58	84.74 \pm 0.84	87.47 \pm 0.50	88.88 \pm 0.33	90.58 \pm 0.18

Table 3: Summary of model additions which did not change or harmed performance. Table entries represent mean \pm std.

Model Variants	10 shot	20 shot	50 shot	100 shot	1K shot
+ maj vote	83.41 \pm 1.23	85.41 \pm 0.88	88.17 \pm 0.53	89.58 \pm 0.57	91.34 \pm 0.34
+ maj vote & idem	83.44 \pm 1.21	85.44 \pm 0.87	88.20 \pm 0.55	89.60 \pm 0.58	91.35 \pm 0.36
+ maj vote all-pairs	83.26 \pm 1.12	85.26 \pm 0.80	88.14 \pm 0.53	89.59 \pm 0.58	91.33 \pm 0.33

Table 4: Summary of alternative majority vote methods. Table entries represent mean \pm std.

6. Conclusion

We have introduced the notion of **latent canonicalization**, where latent representations are manipulated through constrained transformations that set individual factors of variation to fixed values (“canonicalizers”). We show that latent canonicalization encourages representations with markedly better sim2real transfer than comparable models on both the SVHN dataset and on an ImageNet subset, even when only a small sample of the possible combination space was used for training. Notably, latent canonicalized pre-trained models reached few-shot performance with $5\times$ less data than for comparable baselines. Our analysis found that the representation of factors of variation was linearized, as measured by decodability and linear dimensionality reduction (PCA).

We primarily analyzed a realistic but simple SVHN dataset, but also found that latent canonicalization was markedly helpful on an ImageNet subset. The strong performance on both of these datasets (ImageNet in particular) is encouraging for more complex larger-scale data. Interestingly, by focusing the learning on manipulations instead of the samples themselves, we were able to get good performance even with a simplistic simulators that clearly has a non-negligible appearance gap with respect to the real data. This suggests that in the future it would be interesting to try and relax the requirement of a simulator even further.

Our results suggest the promise not only of latent canonicalization, but, more broadly, methods which encourage representational structure by constraining transformations rather than a particular structure itself.

References

- [1] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in beta-vae. *arXiv preprint arXiv:1804.03599*, 2018.
- [2] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [3] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pages 2172–2180, 2016.
- [4] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *Proceedings of*

- the *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
 - [6] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
 - [7] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference in Machine Learning (ICML)*, 2017.
 - [8] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International Conference in Machine Learning (ICML)*, 2015.
 - [9] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations*, 2018.
 - [10] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, Oct 2017.
 - [11] Kaifeng He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2015.
 - [12] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International Conference in Learning Representations (ICLR)*, 2017.
 - [13] Irina Higgins, Arka Pal, Andrei Rusu, Loic Matthey, Christopher Burgess, Alexander Pritzel, Matthew Botvinick, Charles Blundell, and Alexander Lerchner. Darla: Improving zero-shot transfer in reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1480–1490. JMLR. org, 2017.
 - [14] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International Conference in Machine Learning (ICML)*, 2018.
 - [15] Alexander B. Jung, Kentaro Wada, Jon Crall, Satoshi Tanaka, Jake Graving, Sarthak Yadav, Joy Banerjee, Gábor Vecsei, Adam Kraft, Jirka Borovec, Christian Vallentin, Semen Zhydenko, Kilian Pfeiffer, Ben Cook, Ismael Fernández, Weng Chi-Hung, Abner Ayala-Acevedo, Raphael Meudec, Matias Laporte, et al. imgaug. <https://github.com/aleju/imgaug>, 2019. [Online; accessed 14-Sept-2019].
 - [16] Amlan Kar, Aayush Prakash, Ming-Yu Liu, Eric Cameracci, Justin Yuan, Matt Rusiniak, David Acuna, Antonio Torralba, and Sanja Fidler. Meta-sim: Learning to generate synthetic datasets. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019.
 - [17] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018.
 - [18] Minyoung Kim, Yuting Wang, Pritish Sahu, and Vladimir Pavlovic. Bayes-factor-vae: Hierarchical bayesian deep auto-encoder models for factor disentanglement. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2979–2987, 2019.
 - [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 - [20] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *International Conference in Machine Learning (ICML)*, 2015.
 - [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
 - [22] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In *Advances in neural information processing systems*, pages 2539–2547, 2015.
 - [23] Peilun Li, Xiaodan Liang, Daoyuan Jia, and Eric P Xing. Semantic-aware grad-gan for virtual-to-real urban scene adaptation. *arXiv preprint arXiv:1801.01726*, 2018.
 - [24] Yijun Li, Ming-Yu Liu, Xueting Li, Ming-Hsuan Yang, and Jan Kautz. A closed-form solution to photorealistic image stylization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 453–468, 2018.
 - [25] Yuanpeng Li, Liang Zhao, Kenneth Church, and Mohamed Elhoseiny. Compositional continual language learning. 2020.
 - [26] Yuanpeng Li, Liang Zhao, Jianyu Wang, and Joel Hestness. Compositional generalization for primitive substitutions. *arXiv preprint arXiv:1910.02612*, 2019.
 - [27] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Neural Information Processing Symposium (NeurIPS)*, pages 700–708, 2017.
 - [28] Francesco Locatello, Stefan Bauer, Mario Lučić, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Frédéric Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *International Conference on Machine Learning*, 2019. Best Paper Award.
 - [29] Mingshen Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. Learning transferable features with deep adaptation networks. In *International Conference in Machine Learning (ICML)*, 2015.
 - [30] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.
 - [31] Eloi Mehr, Andre Lieutier, Fernando Sanchez Bermudez, Vincent Guitteny, Nicolas Thome, and Matthieu Cord. Manifold learning in quotient spaces. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
 - [32] Franziska Mueller, Florian Bernard, Oleksandr Sotnychenko, Dushyant Mehta, Srinath Sridhar, Dan Casas, and Christian Theobalt. Generated hands for real-time 3d hand tracking from monocular rgb. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 49–59, 2018.
 - [33] Tushar Nagarajan and Kristen Grauman. Attributes as operators: factorizing unseen attribute-object compositions. In *Proceedings of the European Conference on Computer Vision*

- (ECCV), pages 169–185, 2018.
- [34] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
 - [35] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
 - [36] Aayush Prakash, Shaad Boochoon, Mark Brophy, David Acuna, Eric Cameracci, Gavriel State, Omer Shapira, and Stan Birchfield. Structured domain randomization: Bridging the reality gap by context-aware synthetic data, 2018.
 - [37] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
 - [38] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv*, 2018.
 - [39] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European conference on computer vision*, pages 102–118. Springer, 2016.
 - [40] Michal Rolinek, Dominik Zietlow, and Georg Martius. Variational autoencoders pursue pca directions (by accident). *arXiv preprint arXiv:1812.06775*, 2018.
 - [41] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Josh Susskind, Wenda Wang, and Russ Webb. Learning from simulated and unsupervised images through adversarial training. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
 - [42] Jake Snell, Kevin Swerky, and Richard S. Zemel. Prototypical networks for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2017.
 - [43] Martin Sundermeyer, Zoltan-Csaba Marton, Maximilian Durner, Manuel Brucker, and Rudolph Triebel. Implicit 3d orientation learning for 6d object detection from rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 699–715, 2018.
 - [44] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 23–30. IEEE, 2017.
 - [45] Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Adversarial discriminative domain adaptation. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
 - [46] Sjoerd van Steenkiste, Francesco Locatello, Jürgen Schmidhuber, and Olivier Bachem. Are disentangled representations helpful for abstract visual reasoning? *arXiv preprint arXiv:1905.12506*, 2019.
 - [47] Gul Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning from synthetic humans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 109–117, 2017.
 - [48] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, pages 1096–1103, New York, NY, USA, 2008. ACM.
 - [49] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(Dec):3371–3408, 2010.
 - [50] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Neural Information Processing Symposium (NeurIPS)*, 2016.
 - [51] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
 - [52] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.
 - [53] Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Domain adaptation for semantic segmentation via class-balanced self-training. *arXiv preprint arXiv:1810.07911*, 2018.