

PatchVAE: Learning Local Latent Codes for Recognition

Kamal Gupta
University of Maryland College Park
kampta@cs.umd.edu

Saurabh Singh
Google
saurabhsingh@google.com

Abhinav Shrivastava
University of Maryland College Park
abhinav@cs.umd.edu

Abstract

Unsupervised representation learning holds the promise of exploiting large amount of available unlabeled data to learn general representations. A promising technique for unsupervised learning is the framework of Variational Auto-encoders (VAEs). However, unsupervised representations learned by VAEs are significantly outperformed by those learned by supervising for recognition. Our hypothesis is that to learn useful representations for recognition the model needs to be encouraged to learn about repeating and consistent patterns in data. Drawing inspiration from the mid-level representation discovery work, we propose PatchVAE, that reasons about images at patch level. Our key contribution is a bottleneck formulation in a VAE framework that encourages mid-level style representations. Our experiments demonstrate that representations learned by our method perform much better on the recognition tasks compared to those learned by vanilla VAEs.

1. Introduction

Due to the availability of large labeled visual datasets, supervised learning has become the dominant paradigm for visual recognition. That is, to learn about any new concept, the modus operandi is to collect thousands of labeled examples for that concept and train a powerful classifier, such as a deep neural network. This is necessary because the current generation of models based on deep neural networks require large amounts of labeled data [33]. This is in stark contrast to the insights that we have from developmental psychology on how infants develop perception and cognition without any explicit supervision [31]. Moreover, the supervised learning paradigm is ill-suited for applications, such as health care and robotics, where annotated data is hard to obtain either due to privacy concerns or high cost of expert human annotators. In such cases, learning from very few labeled images

or discovering underlying natural patterns in large amounts of unlabeled data can have a large number of potential applications. Discovering such patterns from unlabeled data is the standard setup of unsupervised learning.

Over the past few years, the field of unsupervised learning in computer vision has followed two seemingly different tracks with different goals: generative modeling and self-supervised learning. The goal of generative modeling is to learn the probability distribution from which data was generated, given some training data. A learned model can draw samples from the same distribution or evaluate the likelihoods of new data. Generative models are also useful for learning compact representation of images. However, we argue that these representations are not as useful for visual recognition. This is not surprising since the task of reconstructing images does not require the bottleneck representation to sort out meaningful data useful for recognition and discard the rest; on the contrary, it encourages preserving as much information as possible for reconstruction.

In comparison, the goal in self-supervised learning is to learn representations that are useful for recognition. The standard paradigm is to establish proxy tasks that don't require human-supervision but can provide signals useful for recognition. Due to the mismatch in goals of unsupervised learning for visual recognition and the representations learned from generative modeling, self-supervised learning is a more popular way of learning representations from unlabeled data. However, fundamental limitation of this self-supervised paradigm is that we need to define a proxy-task that can mimic the desired recognition. It is not always possible to establish such a task, nor are these tasks generalizable across recognition tasks.

In this paper, we take the first steps towards enabling the unsupervised generative modeling approach of VAEs to learn representations useful for recognition. Our key hypothesis is that for a representation to be useful, it should capture just the *interesting* parts of the images, as opposed to *everything*

in the images.

What constitutes an interesting image part has been defined and studied in earlier works that pre-date the end-to-end trained deep network methods [30, 7, 14]. Taking inspiration from these works, we propose a novel representation that only encodes such few parts of an image that are repetitive across the dataset, i.e., the patches that occur often in images. By avoiding reconstruction of the entire image our method can focus on regions that are repeating and consistent across many images. In an encoder-decoder based generative model, we constrain the encoder architecture to learn such repetitive parts – both in terms of representations for appearance of these parts (or patches in an image) and where these parts occur. We formulate this using variational auto-encoder (β -VAEs) [19, 23], where we impose novel structure on the latent representations. We use discrete latents to model part presence or absence and continuous latents to model their appearance. We present this approach, PatchVAE, in Section 3 and demonstrate that it learns representations that are much better for recognition as compared to those learned by the standard β -VAEs [19, 23].

In addition, we propose in Section 3.4 that losses that favor foreground, which is more likely to contain repetitive patterns, result in representations that are much better at recognition. In Section 4, we present results on CIFAR100 [20], MIT Indoor Scene Recognition [27], Places [37], and ImageNet [4] datasets. Our contributions are as follows:

- We propose a novel patch-based bottleneck in the VAE framework that learns representations that can encode repetitive parts across images.
- We demonstrate that our method, **PatchVAE**, learns unsupervised representations that are better suited for recognition in comparison to traditional VAEs.
- We show that losses that favor foreground are better for unsupervised learning of representations for recognition.
- We perform extensive ablation analysis to understand the importance of different aspects of the proposed PatchVAE architecture.

2. Related Work

Due to its potential impact, unsupervised learning (particularly for deep networks) is one of the most researched topics in visual recognition over the past few years. Generative models such as VAEs [19, 23, 18, 11], PixelRNN [34], PixelCNN [12, 29], and their variants have proven effective when it comes to learning compressed representation of images while being able to faithfully reconstruct them as well as draw samples from the data distribution. GANs [10, 28, 38, 3] on the other hand, while don't model

the probability density explicitly, can still produce high quality image samples from noise. There has been work combining VAEs and GANs to be able to simultaneously learn image data distribution while being able to generate high quality samples from it [15, 8, 21]. Convolution sparse coding [1] is an alternative approach for reconstruction or image in-painting problems. Our work complements existing generative frameworks in that we provide a structured approach for VAEs that can learn beyond low-level representations. We show the effectiveness of the representations learned by our model by using them for standard visual recognition tasks.

There has been a lot of work in interpreting or disentangling representations learned using generative models such as VAEs [23, 9, 16]. However, there is little evidence of effectiveness of disentangled representations in visual recognition tasks. Semi-supervised Learning using generative models [17, 32], where partial or noisy labels are available to the model during training, has shown lots of promise in applications of generating conditioned samples from the model. In our work however, we focus on incorporating inductive biases in these generative models (e.g., VAEs) such that they can learn representations better suited for visual recognition tasks.

A related, but orthogonal, line of work is self-supervised learning where a proxy task is designed to learn representation useful for recognition. These proxy tasks vary from simple tasks like arranging patches in an image in the correct spatial order [5, 6] and arranging frames from a video in correct temporal order [35, 25], to more involved tasks like in-painting [26] and context prediction [24, 36]. We follow the best practices from this line of work for evaluating the learned representations.

3. Our Approach

Our work builds upon VAE framework proposed by [19]. We briefly review relevant aspects of the VAE framework and then present our approach.

3.1. VAE Review

Standard VAE framework assumes a generative model for data where first a latent \mathbf{z} is sampled from a prior $p(\mathbf{z})$ and then the data is generated from a conditional distribution $G(\mathbf{x}|\mathbf{z})$. A variational approximation $Q(\mathbf{z}|\mathbf{x})$ to the true intractable posterior is introduced and the model is learned by minimizing the following negative variational lower bound (ELBO).

$$\mathcal{L}_{\text{VAE}}(\mathbf{x}) = -\mathbb{E}_{\mathbf{z} \sim Q(\mathbf{z}|\mathbf{x})} [\log G(\mathbf{x}|\mathbf{z})] + \mathbb{KL} [Q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})] \quad (1)$$

$Q(\mathbf{z}|\mathbf{x})$ is often referred to as an encoder as it can be

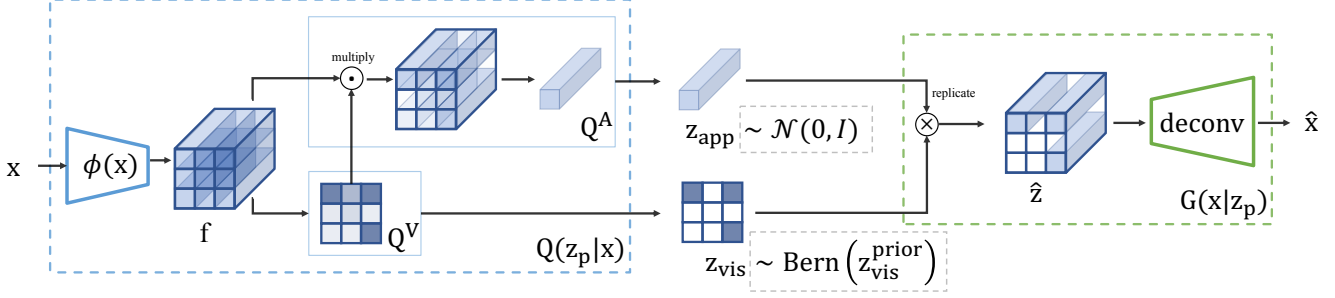


Figure 1: **Proposed PatchVAE Architecture:** Our encoder network computes a set of feature maps \mathbf{f} using $\phi(\mathbf{x})$. This is followed by 2 independent single layer networks - bottom network generates part visibility parameters Q^V . We combine Q^V with output of top network to generate part appearance parameters Q^A . We sample z_{vis} and z_{app} to construct $\hat{\mathbf{z}}$ as described in Section 3.2 which is input to the decoder network. We also visualize the corresponding priors for latents z_{app} and z_{vis} in the dashed gray boxes.

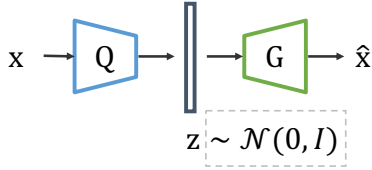


Figure 2: **VAE Architecture:** In a standard VAE architecture, output of encoder network is used to parameterize z . Sample from z is input to the decoder network.

viewed as mapping data to the latent space, while $G(\mathbf{x}|\mathbf{z})$ is referred to as a decoder (or generator) that can be viewed as mapping latents to the data space. Both Q and G are commonly parameterized as neural networks. Fig. 2 shows the commonly used VAE architecture. If the conditional $G(\mathbf{x}|\mathbf{z})$ takes a gaussian form, negative log likelihood in the first term of RHS of Eq. 1 becomes mean squared error between generator output $\hat{\mathbf{x}} = G(\mathbf{x}|\mathbf{z})$ and input data \mathbf{x} . In the second term, prior $p(\mathbf{z})$ is assumed to be a multi-variate normal distribution with zero-mean and diagonal covariance $\mathcal{N}(0, \mathcal{I})$ and the loss simplifies to

$$\mathcal{L}_{\text{VAE}}(\mathbf{x}) = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \mathbb{KL}[Q(\mathbf{z}|\mathbf{x}) \parallel \mathcal{N}(0, \mathcal{I})] \quad (2)$$

When G and Q are differentiable, entire model can be trained with SGD using reparameterization trick [19]. [23] propose an extension for learning disentangled representation by incorporating a weight factor β for the KL Divergence term yielding

$$\mathcal{L}_{\beta\text{VAE}}(\mathbf{x}) = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \beta \mathbb{KL}[Q(\mathbf{z}|\mathbf{x}) \parallel \mathcal{N}(0, \mathcal{I})] \quad (3)$$

VAE framework aims to learn a generative model for the images where the latents \mathbf{z} represent the corresponding low

dimensional generating factors. The latents \mathbf{z} can therefore be treated as image representations that capture the necessary details about images. However, we postulate that representations produced by the standard VAE framework are not ideal for recognition as they are learned to capture *all* details, rather than capturing ‘interesting’ aspects of the data and dropping the rest. This is not surprising since the formulation does not encourage learning semantic information. For learning semantic representations, in the absence of any relevant supervision (as is available in self-supervised approaches), inductive biases have to be introduced. Therefore, taking inspiration from works on unsupervised mid-level pattern discovery [30, 7, 14], we propose a formulation that encourages the encoder to only encode such few parts of an image that are repetitive across the dataset, i.e., the patches that occur often in images.

Since the VAE framework provides a principled way of learning a mapping from image to latent space, we consider it ideal for our proposed extension. We chose β -VAEs for their simplicity and widespread use. In Section 3.2, we describe our approach in detail and in Section 3.4 propose a modification in the reconstruction error computation to bias the error term towards foreground high-energy regions (similar to the biased initial sampling of patterns in [30]).

3.2. PatchVAE

Given an image \mathbf{x} , let $\mathbf{f} = \phi(\mathbf{x})$ be a deterministic mapping that produces a 3D representation \mathbf{f} of size $h \times w \times d_e$, with a total of $L = h \times w$ locations (grid-cells). We aim to encourage the encoder network to only encode parts of an image that correspond to highly repetitive patches. For example, a random patch of noise is unlikely to occur frequently, whereas patterns like faces, wheels, windows, etc. repeat across multiple images. In order capture this intuition, we force the representation \mathbf{f} to be useful for predicting frequently occurring parts in an image, and use *just* these

predicted parts to reconstruct the image. We achieve this by transforming \mathbf{f} to $\hat{\mathbf{z}}$ which encodes a set of parts at a small subset of L locations on the grid cells. We refer to $\hat{\mathbf{z}}$ as “patch latent codes” for an image. Next we describe how we re-tool the β -VAE framework to learn these local latent codes. We first describe our setup for a single part and follow it up with a generalization to multiple parts (Section 3.3).

Image Encoding. Given the image representation $\mathbf{f} = \phi(x)$, we would like to learn part representations at each grid location l (where $l \in \{1, \dots, L\}$). A part is parameterized by its appearance \mathbf{z}_{app} and its visibility $\mathbf{z}_{\text{vis}}^l$ (i.e., presence or absence of the part at grid location l). We use two networks, $Q_{\mathbf{f}}^{\text{A}}$ and $Q_{\mathbf{f}}^{\text{V}}$, to parameterize posterior distributions $Q_{\mathbf{f}}^{\text{A}}(\mathbf{z}_{\text{app}} | \mathbf{f})$ and $Q_{\mathbf{f}}^{\text{V}}(\mathbf{z}_{\text{vis}}^l | \mathbf{f})$ of the part parameters \mathbf{z}_{app} and $\mathbf{z}_{\text{vis}}^l$ respectively. Since the mapping $\mathbf{f} = \phi(\mathbf{x})$ is deterministic, we can re-write these distributions as $Q_{\mathbf{f}}^{\text{A}}(\mathbf{z}_{\text{app}} | \phi(\mathbf{x}))$ and $Q_{\mathbf{f}}^{\text{V}}(\mathbf{z}_{\text{vis}}^l | \phi(\mathbf{x}))$; or simply $Q^{\text{A}}(\mathbf{z}_{\text{app}} | \mathbf{x})$ and $Q^{\text{V}}(\mathbf{z}_{\text{vis}}^l | \mathbf{x})$. Therefore, given an image \mathbf{x} the encoder networks estimate the posterior $Q^{\text{A}}(\mathbf{z}_{\text{app}} | \mathbf{x})$ and $Q^{\text{V}}(\mathbf{z}_{\text{vis}}^l | \mathbf{x})$. Note that \mathbf{f} is a deterministic feature map, whereas \mathbf{z}_{app} and $\mathbf{z}_{\text{vis}}^l$ are stochastic.

Image Decoding. We utilize a generator or decoder network G , that given \mathbf{z}_{vis} and \mathbf{z}_{app} , reconstructs the image. First, we sample a part appearance $\hat{\mathbf{z}}_{\text{app}}$ (d_p dimensional, continuous) and then sample part visibilities $\hat{\mathbf{z}}_{\text{vis}}^l$ (L dimensional, binary) one for each location l from the posteriors

$$\begin{aligned}\hat{\mathbf{z}}_{\text{app}} &\sim Q^{\text{A}}(\mathbf{z}_{\text{app}} | \mathbf{x}) \\ \hat{\mathbf{z}}_{\text{vis}}^l &\sim Q^{\text{V}}(\mathbf{z}_{\text{vis}}^l | \mathbf{x}), \quad \text{where } l \in \{1, \dots, L\}\end{aligned}\quad (4)$$

Next, we construct a 3D representation $\hat{\mathbf{z}}$ by placing $\hat{\mathbf{z}}_{\text{app}}$ at every location l where the part is present (i.e., $\hat{\mathbf{z}}_{\text{vis}}^l = 1$). This can be implemented by a broadcasted product of $\hat{\mathbf{z}}_{\text{app}}$ and $\hat{\mathbf{z}}_{\text{vis}}^l$. We refer to $\hat{\mathbf{z}}$ as **patch latent code**. Again note that \mathbf{f} is deterministic and $\hat{\mathbf{z}}$ is stochastic. Finally, a deconvolutional network takes $\hat{\mathbf{z}}$ as input and generates an image $\hat{\mathbf{x}}$. This image generation process can be written as

$$\hat{\mathbf{x}} \sim G(\mathbf{x} | \mathbf{z}_{\text{vis}}^1, \mathbf{z}_{\text{vis}}^2, \dots, \mathbf{z}_{\text{vis}}^L, \mathbf{z}_{\text{app}}) \quad (5)$$

Since all latent variables ($\mathbf{z}_{\text{vis}}^l$ for all l and \mathbf{z}_{app}) are independent of each other, they can be stacked as

$$\mathbf{z}_{\text{p}} = [\mathbf{z}_{\text{vis}}^1; \mathbf{z}_{\text{vis}}^2; \dots; \mathbf{z}_{\text{vis}}^L; \mathbf{z}_{\text{app}}]. \quad (6)$$

This enables us to use a simplified notation (refer to (4) and (5)):

$$\begin{aligned}\hat{\mathbf{z}}_{\text{p}} &\sim Q^{\{\text{A}, \text{V}\}}(\mathbf{z}_{\text{p}} | \mathbf{x}) \\ \hat{\mathbf{x}} &\sim G(\mathbf{x} | \mathbf{z}_{\text{p}})\end{aligned}\quad (7)$$

Note that despite the additional structure, our model still resembles the setup of variational auto-encoders. The primary difference arises from: (1) use of discrete latents for part visibility, (2) patch-based bottleneck imposing additional structure on latents, and (4) feature assembly for generator.

Training. We use the training setup of β -VAE and use the maximization of variational lower bound to train the encoder and decoder jointly (described in Section 3.1). The posterior Q^{A} , which captures the appearance of a part, is assumed to be a zero-mean Normal distribution with diagonal covariance $\mathcal{N}(0, \mathcal{I})$. The posterior Q^{V} , which captures the presence or absence a part, is assumed to be a Bernoulli distribution $\text{Bern}(\mathbf{z}_{\text{vis}}^{\text{prior}})$ with prior $\mathbf{z}_{\text{vis}}^{\text{prior}}$. Therefore, the ELBO for our approach can be written as (refer to (3)):

$$\begin{aligned}\mathcal{L}_{\text{PatchVAE}}(\mathbf{x}) &= -\mathbb{E}_{\mathbf{z}_{\text{p}} \sim Q^{\{\text{A}, \text{V}\}}(\mathbf{z}_{\text{p}} | \mathbf{x})} [G(\mathbf{x} | \mathbf{z}_{\text{p}})] \\ &\quad + \beta \mathbb{KL} [Q^{\{\text{A}, \text{V}\}}(\mathbf{z}_{\text{p}} | \mathbf{x}) \parallel p(\mathbf{z}_{\text{p}})]\end{aligned}\quad (8)$$

where, the \mathbb{KL} term can be expanded as:

$$\begin{aligned}\mathbb{KL} [Q^{\{\text{A}, \text{V}\}}(\mathbf{z}_{\text{p}} | \mathbf{x}) \parallel p(\mathbf{z}_{\text{p}})] &= \\ \beta_{\text{app}} \sum_{l=1}^L \mathbb{KL} (Q^{\text{V}}(\mathbf{z}_{\text{vis}}^l | \mathbf{x}) \parallel \text{Bern}(\mathbf{z}_{\text{vis}}^{\text{prior}})) &\quad (9) \\ + \beta_{\text{vis}} \mathbb{KL} (Q^{\text{A}}(\mathbf{z}_{\text{app}} | \mathbf{x}) \parallel \mathcal{N}(0, \mathcal{I}))\end{aligned}$$

Implementation details. As discussed in Section 3.1, the first and second terms of the RHS of (8) can be trained using L2 reconstruction loss and reparameterization trick [19]. In addition, we also need to compute KL Divergence loss for part visibility. Learning discrete probability distribution is a challenging task since there is no gradient defined to back-propagate reconstruction loss through the stochastic layer at decoder even when using the reparameterization trick. Therefore, we use the relaxed-bernoulli approximation [22, 2] for training part visibility distributions $\mathbf{z}_{\text{vis}}^l$.

For an $H \times W$ image, network $Q(\mathbf{f} | \mathbf{x})$ first generates feature maps of size $(h \times w \times d_e)$, where (h, w) are spatial dimensions and d_e is the number of channels. Therefore, the number of locations $L = h \times w$. Encoders $Q_{\mathbf{f}}^{\text{A}}(\mathbf{z}_{\text{app}} | \mathbf{f})$ and $Q_{\mathbf{f}}^{\text{V}}(\mathbf{z}_{\text{vis}}^l | \mathbf{f})$ are single layer neural networks to compute \mathbf{z}_{app} and $\mathbf{z}_{\text{vis}}^l$. $\mathbf{z}_{\text{vis}}^l$ is $(h \times w \times 1)$ -dimensional multivariate bernoulli parameter and \mathbf{z}_{app} is $(1 \times 1 \times d_p)$ -dimensional multivariate gaussian. d_p is length of the latent vector for a single part. Input to the decoder $\hat{\mathbf{z}}$ is $(h \times w \times d_p)$ -dimensional. In all our experiments, we fix $h = \frac{H}{8}$ and $w = \frac{W}{8}$.

Constructing \mathbf{z}_{app} . Notice that \mathbf{f} is an $(h \times w \times d_e)$ -dimensional feature map and $\mathbf{z}_{\text{vis}}^l$ is $(h \times w \times 1)$ -

dimensional binary output, but \mathbf{z}_{app} is $(1 \times 1 \times d_p)$ -dimensional feature vector. If $\sum_l \mathbf{z}_{\text{vis}}^l > 1$, the part occurs at multiple locations in an image. Since all these locations correspond to same part, their appearance should be the same. To incorporate this, we take the weighted average of the part appearance feature at each location, weighted by the probability that the part is present. Since we use the probability values for averaging the result is deterministic. This operation is encapsulated by the Q^A encoder (refer to Figure 1). During image generation, we sample $\hat{\mathbf{z}}_{\text{app}}$ once and replicate it at each location where $\hat{\mathbf{z}}_{\text{vis}}^l = 1$. During training, this forces the model to: (1) only predict $\hat{\mathbf{z}}_{\text{vis}}^l = 1$ where similar looking parts occur, and (2) learn a common representation for the part that occurs at these locations. Note that \mathbf{z}_{app} can be modeled as a mixture of distributions (e.g., mixture of gaussians) to capture complicated appearances. However, in this work we assume that the convolutional neural network based encoders are powerful enough to map variable appearance of semantic concepts to similar feature representations. Therefore, we restrict ourselves to a single gaussian distribution.

3.3. PatchVAE with multiple parts

Next we extend the framework described above to use multiple parts. To use N parts, we use $N \times 2$ encoder networks $Q^{A(i)}(\mathbf{z}_{\text{app}}^{(i)} | \mathbf{x})$ and $Q^{V(i)}(\mathbf{z}_{\text{vis}}^{(i)} | \mathbf{x})$, where $\mathbf{z}_{\text{app}}^{(i)}$ and $\mathbf{z}_{\text{vis}}^{(i)}$ parameterize the i^{th} part. Again, this can be implemented efficiently as 2 networks by concatenating the outputs together. The image generator samples $\hat{\mathbf{z}}_{\text{app}}^{(i)}$ and $\hat{\mathbf{z}}_{\text{vis}}^{(i)}$ from the outputs of these encoder networks and constructs $\hat{\mathbf{z}}^{(i)}$. We obtain the final **patch latent code** $\hat{\mathbf{z}}$ by concatenating all $\hat{\mathbf{z}}^{(i)}$ in channel dimension. Therefore, $\hat{\mathbf{z}}^{(i)}$ is $(h \times w \times d_p)$ -dimensional and $\hat{\mathbf{z}}$ is $(h \times w \times (N \times d_p))$ -dimensional stochastic feature map. For this multiple part case, (6) can be written as:

$$\mathbf{z}_p = [\mathbf{z}_p^{(1)}; \mathbf{z}_p^{(1)}; \dots; \mathbf{z}_p^{(N)}] \quad (10)$$

where $\mathbf{z}_p^{(i)} = [\mathbf{z}_{\text{vis}}^{1(i)}; \mathbf{z}_{\text{vis}}^{1(i)}; \dots; \mathbf{z}_{\text{vis}}^{L(i)}; \mathbf{z}_{\text{app}}^{(i)}]$.

Similarly, (8) and (9) can be written as:

$$\begin{aligned} \mathcal{L}_{\text{MultiPatchVAE}}(\mathbf{x}) = & -\mathbb{E}_{\mathbf{z}_p} [G(\mathbf{x} | \mathbf{z}_p)] \\ & + \beta_{\text{app}} \sum_{i=1}^N \sum_{l=1}^L \mathbb{KL} \left(Q^{V(i)}(\mathbf{z}_{\text{vis}}^{l(i)} | \mathbf{x}) \parallel \text{Bern}(\mathbf{z}_{\text{vis}}^{\text{prior}}) \right) \\ & + \beta_{\text{vis}} \sum_{i=1}^N \mathbb{KL} \left(Q^{A(i)}(\mathbf{z}_{\text{app}}^{(i)} | \mathbf{x}) \parallel \mathcal{N}(0, \mathcal{I}) \right) \end{aligned} \quad (11)$$

The training details and assumptions of posteriors follow the previous section.

3.4. Improved Reconstruction Loss

The L2 reconstruction loss used for training β -VAEs (and other reconstruction based approaches) gives equal importance to each region of an image. This might be reasonable for tasks like image compression and image de-noising. However, for the purposes of learning semantic representations, not all regions are equally important. For example, “sky” and “walls” occupy large portions of an image, whereas concepts like “windows,” “wheels,” “faces” are comparatively smaller, but arguably more important. To incorporate this intuition, we use a simple and intuitive strategy to weigh the regions in an image in proportion to the gradient energy in the region. More concretely, we compute laplacian of an image to get the intensity of gradients per-pixel and average the gradient magnitudes in 8×8 local patches. The weight multiplier for the reconstruction loss of each 8×8 patch in the image is proportional to the average magnitude of the patch. All weights are normalized to sum to one. We refer to this as **weighted loss** (\mathcal{L}_w). Note that this is similar to the gradient-energy biased sampling of mid-level patches used in [30, 7]. We show examples of weight masks for some of the images in the supplementary material.

In addition, we also consider an adversarial training strategy from GANs to train VAEs as proposed by [21], where the discriminator network from GAN implicitly learns to compare images and gives a more abstract reconstruction error for the VAE. We refer to this variant by using ‘GAN’ suffix in experiments. In Section 4, we demonstrate that the proposed weighted loss (\mathcal{L}_w) is complementary to the discriminator loss from adversarial training, and these losses result in better recognition capabilities for both β -VAE and PatchVAE.

4. Experiments

Datasets. We evaluate our proposed model on CIFAR100 [20], MIT Indoor Scene Recognition [27], Places [37] and Imagenet [4] datasets. CIFAR100 consists of 60000 32×32 color images in 100 classes, with 600 images per class. There are 50000 training images and 10000 test images. Indoor dataset contains 67 categories, and a total of 15620 images. Train and test subsets consist of 80 and 20 images per class respectively. Places dataset has 2.5 millions of images with 205 categories. Imagenet dataset has over a million images from 1000 categories.

Learning paradigm. In order to evaluate the utility of features learned for recognition, we setup the learning paradigm as follows: we will first train the model in an unsupervised manner on all the images other than test set images. After that, we discard the generator network and use only part of the encoder network $\phi(\mathbf{x})$ to train a supervised model on the classification task of the respective dataset. We study

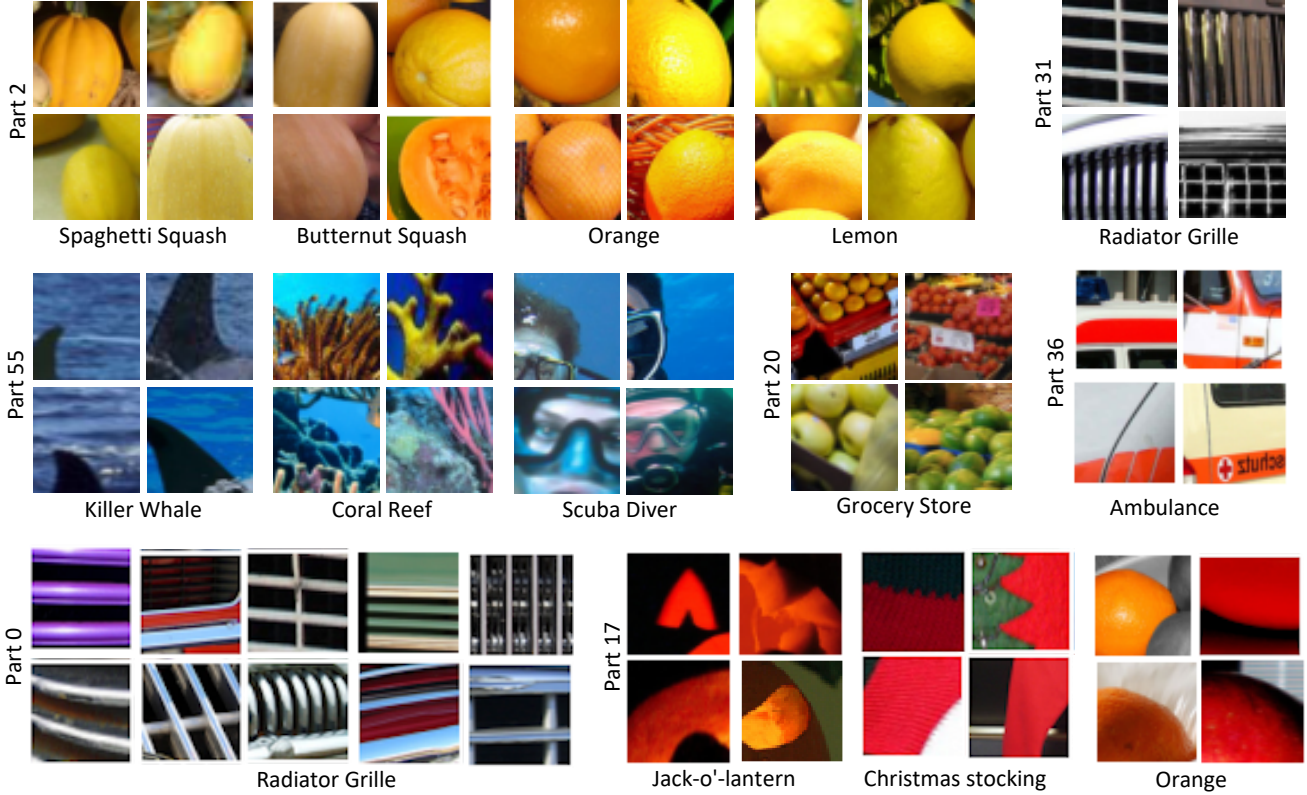


Figure 3: Concepts captured by parts: We visualize a few representative examples for several parts to qualitatively demonstrate the visual concepts captured by parts. For each part, we crop image patches centered on the part location where it is predicted to be present. Selected patches are sorted by part visibility probability as score. We have manually selected a diverse set from the top 50 occurrences from the training images. As visible, a single part may capture diverse set of concepts that are similar in shape or texture or occur in similar context, but belong to different categories. We show which categories the patches come from (note that category information was not used while training the model).

different training strategies for the classification stage as discussed later.

Training details. In all experiments, we use the following architectures. For CIFAR100, Indoor67, and Place205, $\phi(\mathbf{x})$ has a conv layer followed by two residual blocks [13]. For ImageNet, $\phi(\mathbf{x})$ is a ResNet18 model (a conv layer followed by four residual blocks). For all datasets, Q^A and Q^V have a single conv layer each. For classification, we start from $\phi(\mathbf{x})$, and add a fully-connected layer with 512 hidden units and a final fully-connected layer as classifier. More details can be found in the supplementary material

During the unsupervised learning part of training, all methods are trained for 90 epochs for CIFAR100 and Indoor67, 2 epochs for Places205, and 30 epochs for ImageNet dataset. All methods use ADAM optimizer for training, with initial learning rate of 1×10^{-4} and a minibatch size of 128. For relaxed bernoulli in Q^V , we start with the temperature of 1.0 with an annealing rate of 3×10^{-5} (details in [2]). For

training the classifier, all methods use stochastic gradient descent (SGD) with momentum with a minibatch size of 128. Initial learning rate is 1×10^{-2} and we reduce it by a factor of 10 every 30 epochs. All experiments are trained for 90 epochs for CIFAR100 and Indoor67, 5 epochs for Places205, and 30 epochs for ImageNet datasets.

Baselines. We use the β -VAE model (Section 3.1) as our primary baseline. In addition, we use weighted loss and discriminator loss resulting in the β -VAE-* family of baselines. We also compare against a BiGAN model from [8]. We use similar backbone architectures for encoder/decoder (and discriminator if present) across all methods, and tried to keep the number of parameters in different approaches comparable to the best of our ability. Exact architecture details can be found in the supplementary material.

4.1. Results

In Table 1, we report the top-1 classification results on CIFAR100, Indoor67, and Places205 datasets for all meth-

Table 1: Classification results on CIFAR100, Indoor67, and Places205. We initialize the classification model with the representations $\phi(\mathbf{x})$ learned from unsupervised learning task. The model $\phi(\mathbf{x})$ comprises of a conv layer followed by two residual blocks (each having 2 conv layers). First column (called ‘Conv1’) corresponds to Top-1 classification accuracy with pre-trained model with the first conv layer frozen, second and third columns correspond to results with 3 conv layers and 5 conv layers frozen respectively. Details in Section 4.1.

Model	CIFAR100			Indoor67			Places205		
	Conv1	Conv3	Conv5	Conv1	Conv3	Conv5	Conv1	Conv3	Conv5
β -VAE	44.12	39.65	28.57	20.08	17.76	13.06	28.29	24.34	8.89
β -VAE + \mathcal{L}_w	44.96	40.30	28.33	21.34	19.48	13.96	29.43	24.93	9.41
β -VAE-GAN	44.69	40.13	29.89	19.10	17.84	13.06	28.48	24.51	9.72
β -VAE-GAN + \mathcal{L}_w	45.61	41.35	31.53	20.45	18.36	14.33	29.63	25.26	10.66
PatchVAE	43.07	38.58	28.72	20.97	19.18	13.43	28.63	24.95	11.09
PatchVAE + \mathcal{L}_w	43.75	40.37	30.55	23.21	21.87	15.45	29.39	26.29	12.07
PatchVAE-GAN	44.45	40.57	31.74	21.12	19.63	14.55	28.87	25.25	12.21
PatchVAE-GAN + \mathcal{L}_w	45.39	41.74	32.65	22.46	21.87	16.42	29.36	26.30	13.39
BiGAN	47.72	41.89	31.58	21.64	17.09	9.70	30.06	25.11	10.82
Imagenet Pretrained	55.99	54.99	54.36	45.90	45.82	40.90	37.08	36.46	31.26

Table 2: ImageNet classification results using ResNet18. We initialize weights from using the unsupervised task and fine-tune the last 2 residual blocks. Details in Section 4.1.

Model	Top-1 Acc.	Top-5 Acc.
β -VAE	44.45	69.67
PatchVAE	47.01	71.71
β -VAE + \mathcal{L}_w	47.28	71.78
PatchVAE + \mathcal{L}_w	47.87	72.49
Imagenet Supervised	61.37	83.79

ods with different training strategies for classification. First, we keep all the pre-trained weights in $\phi(\mathbf{x})$ from the unsupervised task frozen and only train the two newly added conv layers in the classification network (reported under column ‘Conv5’). We notice that our method (with different losses) generally outperforms the β -VAE counterpart by a healthy margin. This shows that the representations learned by PatchVAE framework are better for recognition compared to β -VAEs. Moreover, better reconstruction losses (‘GAN’ and \mathcal{L}_w) generally improve both β -VAE and PatchVAE, and are complementary to each other.

Next, we fine-tune the last residual block along with the two conv layers (‘Conv3’ column). We observe that PatchVAE performs better than VAE under all settings except

the for CIFAR100 with just L2 loss. However, when using better reconstruction losses, the performance of PatchVAE improves over β -VAE. Similarly, we fine-tune all but the first conv layer and report the results in ‘Conv1’ column. Again, we notice similar trends, where our method generally performs better than β -VAE on Indoor67 and Places205 dataset, but β -VAE performs better CIFAR100 by a small margin. When compared to BiGAN, PatchVAE representations are better on all datasets (‘Conv5’) by a huge margin. However, when fine-tuning the pre-trained weights, BiGAN performs better on two out of four datasets. We also report results using pre-trained weights in $\phi(\mathbf{x})$ using *supervised* ImageNet classification task (last column, Table 1) for completeness. The results indicate that PatchVAE learns better semantic representations compared to β -VAE.

ImageNet Results. Finally, we report results on the large-scale ImageNet benchmark in Table 2. For these experiments, we use ResNet18 [13] architecture for all methods. All weights are first learned using the unsupervised tasks. Then, we fine-tune the last two residual blocks and train the two newly added conv layers in the classification network (therefore, first conv layer and the following two residual blocks are frozen). We notice that PatchVAE framework outperforms β -VAE under all settings, and the proposed weighted loss helps both approaches. Finally, the last row in Table 2 reports classification results of same architecture randomly initialized and trained end-to-end on ImageNet

Table 3: Effect of maximum number of patches (N): Increasing N increases the discriminative power in case of CIFAR100 but has little or negative effect in case of Indoor67

N	Acc on CIFAR100	Acc on Indoor67
4	27.59	14.40
8	28.74	12.69
16	28.94	14.33
32	27.78	13.28
64	29.00	12.76

Table 4: Effect of d_p : Increasing d_p or the number of hidden units for a patch has very little impact on classification performance of base classifier.

d_p	Acc on CIFAR100	Acc on Indoor67
3	28.63	14.25
6	28.97	14.55
9	28.21	14.55

using supervised training for comparison.

4.2. Ablation Studies

We study the impact of various hyper-parameters used in our experiments. For the purpose of this evaluation, we follow a similar approach as in the ‘Conv5’ column of Table 1 and all hyperparameters from the previous section. We use CIFAR100 and Indoor67 datasets for ablation analysis.

Maximum number of patches. Maximum number of parts N used in our framework. Depending on the dataset, higher value of N can provide wider pool of patches to pick from. However, it can also make the unsupervised learning task harder, since in a minibatch of images, we might not get too many repeat patches. Table 3(left) shows the effect of N on CIFAR100 and Indoor67 datasets. We observe that while increasing number of patches improves the discriminative power in case of CIFAR100, it has little or negative effect in case of Indoor67. A possible reason for this decline in performance for Indoor67 can be smaller size of the dataset (i.e., fewer images to learn).

Number of hidden units for a patch appearance \hat{z}_{app} . Next, we study the impact of the number of channels in the appearance feature \hat{z}_{app} for each patch (d_p). This parameter reflects the capacity of individual patch’s latent representation. While this parameter impacts the reconstruction quality of images. We observed that it has little or no effect on

Table 5: Effect of z_{vis}^{prior} : Increasing z_{vis}^{prior} or the prior on patch visibility has adverse effect on classification performance.

z_{vis}^{prior}	Acc on CIFAR100	Acc on Indoor67
0.01	28.86	14.33
0.05	28.67	14.25
0.1	28.31	14.03

Table 6: Effect of β_{vis} : Both high and low values of β_{vis} can deteriorate the performance of learned representations on classification.

β_{vis}	Acc on CIFAR100	Acc on Indoor67
0.06	30.11	14.10
0.3	30.37	15.67
0.6	28.90	13.51

the classification performance of the base features. Results are summarized in Table 4(right) for both CIFAR100 and Indoor67 datasets.

Prior probability for patch visibility z_{vis}^{prior} . In all our experiments, prior probability for a patch is fixed to $1/N$, i.e., inverse of maximum number of patches. The intuition is to encourage each location on visibility maps to fire for at most one patch. Increasing this patch visibility prior will allow all patches to fire at the same location. While this would make the reconstruction task easier, it will become harder for individual patches to capture anything meaningful. Table 5 shows the deterioration of classification performance on increasing z_{vis}^{prior} .

Patch visibility loss weight β_{vis} . The weight for patch visibility KL Divergence has to be chosen carefully. If β_{vis} is too low, more patches can fire at same location and this harms the the learning capability of patches; and if β_{vis} is too high, decoder will not receive any patches to reconstruct from and both reconstruction and classification will suffer. Table 6 summarizes the impact of varying β_{vis} .

5. Conclusion

We presented a patch-based bottleneck in a VAE framework that encourages learning useful representations for recognition. Our method, PatchVAE, constrains the encoder architecture to only learn patches that are repetitive and consistent in images as opposed to learning *everything*, and therefore results in representations that perform much better for recognition tasks compared to vanilla VAEs. We also

demonstrate that losses that favor high-energy foreground regions of an image are better for unsupervised learning of representations for recognition.

References

- [1] Lama Affara, Bernard Ghanem, and Peter Wonka. Supervised convolutional sparse coding. *CoRR*, abs/1804.02678, 2018. 2
- [2] Eirikur Agustsson, Fabian Mentzer, Michael Tschannen, Lukas Cavigelli, Radu Timofte, Luca Benini, and Luc V Gool. Soft-to-hard vector quantization for end-to-end learning compressible representations. In *Advances in Neural Information Processing Systems*, pages 1141–1151, 2017. 4, 6
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017. 2
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. IEEE Conference on*, pages 248–255. Ieee, 2009. 2, 5
- [5] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Context as supervisory signal: Discovering objects with predictable context. In *European Conference on Computer Vision*, pages 362–377. Springer, 2014. 2
- [6] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. In *International Conference on Computer Vision (ICCV)*, 2015. 2
- [7] Carl Doersch, Saurabh Singh, Abhinav Gupta, Josef Sivic, and Alexei A. Efros. What makes paris look like paris? *ACM Transactions on Graphics (SIGGRAPH)*, 31(4):101:1–101:9, 2012. 2, 3, 5
- [8] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016. 2, 6
- [9] Marco Fraccaro, Simon Kamronn, Ulrich Paquet, and Ole Winther. A disentangled recognition and nonlinear dynamics model for unsupervised learning. In *NIPS*, 2017. 2
- [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2
- [11] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra. Draw: A recurrent neural network for image generation. *arXiv preprint arXiv:1502.04623*, 2015. 2
- [12] Ishaan Gulrajani, Kundan Kumar, Faruk Ahmed, Adrien Ali Taiga, Francesco Visin, David Vazquez, and Aaron Courville. Pixelvae: A latent variable model for natural images. *arXiv preprint arXiv:1611.05013*, 2016. 2
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6, 7
- [14] Mayank Juneja, Andrea Vedaldi, CV Jawahar, and Andrew Zisserman. Blocks that shout: Distinctive parts for scene classification. In *Computer Vision and Pattern Recognition*, pages 923–930, 2013. 2, 3
- [15] Salman H Khan, Munawar Hayat, and Nick Barnes. Adversarial training of variational auto-encoders for high fidelity image generation. *arXiv preprint arXiv:1804.10323*, 2018. 2
- [16] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018. 2
- [17] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, pages 3581–3589, 2014. 2
- [18] Diederik P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, pages 4743–4751, 2016. 2
- [19] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2, 3, 4
- [20] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009. 2, 5
- [21] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300*, 2015. 2, 5
- [22] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016. 4
- [23] Loïc 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 *ICLR 2017*, 2017. 2, 3
- [24] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016. 2
- [25] Deepak Pathak, Ross Girshick, Piotr Dollár, Trevor Darrell, and Bharath Hariharan. Learning features by watching objects move. In *CVPR*, 2017. 2
- [26] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2536–2544, 2016. 2
- [27] Ariadna Quattoni and Antonio Torralba. Recognizing indoor scenes. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 413–420. IEEE, 2009. 2, 5
- [28] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015. 2
- [29] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P. Kingma. Pixelcnn++: A pixelcnn implementation with discretized logistic mixture likelihood and other modifications. In *ICLR*, 2017. 2

- [30] Saurabh Singh, Abhinav Gupta, and Alexei A. Efros. Unsupervised discovery of mid-level discriminative patches. In *European Conference on Computer Vision*, 2012. 2, 3, 5
- [31] Linda Smith and Michael Gasser. The development of embodied cognition: Six lessons from babies. *Artificial life*, 11(1-2):13–29, 2005. 1
- [32] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in neural information processing systems*, pages 3483–3491, 2015. 2
- [33] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 843–852, 2017. 1
- [34] Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. *CoRR*, abs/1601.06759, 2016. 2
- [35] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2015. 2
- [36] Xiaolong Wang, Kaiming He, and Abhinav Gupta. Transitive invariance for self-supervised visual representation learning. In *Proc. of Int’l Conf. on Computer Vision (ICCV)*, 2017. 2
- [37] 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*, 40(6):1452–1464, 2017. 2, 5
- [38] 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. 2