

# Stacked Generative Adversarial Networks

Xun Huang<sup>1</sup> Yixuan Li<sup>2</sup> Omid Poursaeed<sup>2</sup> John Hopcroft<sup>1</sup> Serge Belongie<sup>1,3</sup>

<sup>1</sup>Department of Computer Science, Cornell University

<sup>2</sup>School of Electrical and Computer Engineering, Cornell University

<sup>3</sup>Cornell Tech

{xh258, yl2363, op63, sjb344}@cornell.edu jeh@cs.cornell.edu

## Abstract

In this paper we aim to leverage the powerful bottom-up discriminative representations to guide a top-down generative model. We propose a novel generative model named *Stacked Generative Adversarial Networks* (SGAN), which is trained to invert the hierarchical representations of a discriminative bottom-up deep network. Our model consists of a top-down stack of GANs, each trained to generate “plausible” lower-level representations, conditioned on higher-level representations. A representation discriminator is introduced at each feature hierarchy to encourage the representation manifold of the generator to align with that of the bottom-up discriminative network, providing intermediate supervision. In addition, we introduce a conditional loss that encourages the use of conditional information from the layer above, and a novel entropy loss that maximizes a variational lower bound on the conditional entropy of generator outputs. To the best of our knowledge, the entropy loss is the first attempt to tackle the conditional model collapse problem that is common in conditional GANs. We first train each GAN of the stack independently, and then we train the stack end-to-end. Unlike the original GAN that uses a single noise vector to represent all the variations, our SGAN decomposes variations into multiple levels and gradually resolves uncertainties in the top-down generative process. Experiments demonstrate that SGAN is able to generate diverse and high-quality images, as well as being more interpretable than a vanilla GAN.

## 1. Introduction

Recent years have witnessed tremendous success of deep neural networks (DNNs), especially the kind of bottom-up neural networks that are trained for discriminative tasks. In particular, Convolutional Neural Networks (CNNs) have achieved impressive accuracy on the challenging ImageNet classification benchmark [30, 58, 59, 21, 53]. Interestingly, it has been shown that CNNs trained on ImageNet for classification can learn representations that are transferable to

other tasks [56], and even to other modalities [20]. However, bottom-up discriminative models are focused on learning useful representations from data, being incapable of capturing the data distribution.

Learning top-down generative models that can explain complex data distribution is a long-standing problem in machine learning research. The expressive power of deep neural networks makes them natural candidates for generative models, and several recent works have shown promising results [28, 17, 45, 36, 70, 38, 9]. While state-of-the-art DNNs can rival human performance in discriminative tasks such as object classification, current best deep generative models still fail when there are large variations in the data distribution.

A natural question therefore arises: can we leverage the hierarchical representations in a discriminatively trained model to help the learning of top-down generative models? In this paper, we propose a novel generative model named *Stacked Generative Adversarial Networks* (SGAN). Our model consists of a top-down stack of GANs, each trained to generate “plausible” lower-level representations conditioned on higher-level representations. Similar to the image discriminator in the original GAN model which is trained to distinguish “fake” images from “real” ones, we introduce a set of *representation discriminators* that are trained to distinguish “fake” representations from “real” representations. The *adversarial loss* introduced by the representation discriminator forces the intermediate representations of the SGAN to lie on the manifold of the bottom-up DNN’s representation space. In addition to the adversarial loss, we also introduce a *conditional loss* that imposes each generator to use the conditional information from the layer above, and a novel *entropy loss* that encourages each generator to generate diverse representations by maximizing a lower bound for the conditional entropy of the generator outputs. By stacking several GANs in a top-down way and using the top-most GAN to receive labels and the bottom-most GAN to generate images, SGAN can be trained to model the data distribution conditioned on class labels. Compared with a vanilla GAN, SGAN has three sub-

stantial advantages:

1. **Quality.** SGAN decomposes the difficult image generation task into several easier sub-tasks by leveraging the discriminative representations. As a result, SGAN can generate higher-quality samples than a vanilla GAN. Empirically, we also observe SGAN to be more stable than a vanilla GAN.
2. **Diversity.** With the proposed entropy loss, SGAN can avoid the (conditional) *model collapse* phenomenon that is common in (conditional) GAN training, thus being able to generate more diverse samples.
3. **Interpretability.** Our SGAN is more interpretable than a vanilla GAN in that it uses different levels of noise variables to represent different levels of variations.

## 2. Related Work

**Deep Generative Image Models.** There has been a large body of work on generative image modeling with deep learning. Some early efforts include Restricted Boltzmann Machines (RBMs) [22] and Deep Belief Networks (DBNs) [23]. More recently, several successful paradigms of deep generative models have emerged, including the auto-regressive models [32, 16, 60, 45, 46, 19], Variational Auto-encoders (VAEs) [28, 27, 51, 66, 18], and Generative Adversarial Networks (GANs) [17, 5, 48, 50, 54, 33]. Our work builds upon the GAN framework, which uses a generator to transform a noise vector into an image and a discriminator to distinguish between “real” and “generated” images. The generator is trained to generate images that are real enough to “fool” the discriminator.

However, due to the vast variations in image content, it is still challenging for GANs to generate diverse images with sufficient details. To this end, several works have attempted to factorize a GAN into a series of GANs, decomposing the difficult task into several more tractable sub-tasks. [5] proposes a LAPGAN model that factorizes the generative process into multi-resolution GANs, with each GAN generating a higher-resolution residual conditioned on a lower-resolution image. Although both LAPGAN and SGAN consist of a sequence of GANs each working at one scale, LAPGAN focuses on generating *multi-resolution images* from coarse to fine while our SGAN aims at modeling *multi-level representations* from high-level to low-level. [64] proposes a S<sup>2</sup>-GAN model, using one GAN to generate surface normals and another GAN to generate images conditioned on surface normals. Surface normals can be viewed as a specific type of image representations, capturing the underlying 3D structure of an indoor scene. However, our framework can leverage the more general and powerful multi-level representations in a pre-trained discriminative DNN.

There are several works that use a pre-trained discriminative model to aid the training of a generative model. [31, 7] add a discriminative regularization term that encourages the image reconstructed by a VAE to be similar to the original image in the representation space defined by a pre-trained bottom-up DNN. [14, 13, 61, 26] use an additional “style loss” which minimizes the *L2* distance between Gram matrices on some representation space. All the methods above only add loss terms to regularize the *output* of the generator, without regularizing its *internal representations*. Our SGAN adopts a very different approach, using adversarial training to encourage the internal representations of the generator to reside on the representation manifold defined by a pre-trained discriminative model.

**Matching Intermediate Representations Between Two DNNs.** There have been research directions attempting to “match” the intermediate representations between two DNNs. [52, 20] use the intermediate representations of one pre-trained DNN to guide another DNN in the context of knowledge transfer. Our method can be considered as a special kind of knowledge transfer. However, we aim at transferring the knowledge in a bottom-up DNN to a top-down generative model, instead of another bottom-up DNN. Another direction is the layer-wise reconstruction loss used in some auto-encoder architectures [62, 49, 69, 68]. The layer-wise loss is usually accompanied by lateral connections from the encoder to the decoder at each level of the hierarchy. Our SGAN, however, is a fully independent generative model and does not require any information from the encoder once training completes. Another important difference is that we use adversarial loss instead of *L2* reconstruction loss to match intermediate representations.

**Visualizing Deep Representations.** Our work is also related to the recent efforts in visualizing and understanding the internal representations of DNNs. One popular approach uses gradient-based optimization to find an image whose representation is close to the representation we want to visualize [57, 37, 43]. Other approaches, such as [8], train a top-down deconvolutional network to reconstruct the input image from a feature representation by minimizing the Euclidean reconstruction error in image space. However, there is inherent uncertainty in the reconstruction process, since the representations in higher layers of the DNN are trained to be invariant to irrelevant transformations and to ignore low-level details. With Euclidean training objective, the deconvolutional network learns to produce a blurry image when it is uncertain about the details.

To alleviate this problem, [7] further proposes a feature loss that encourages the reconstructed image to have similar high-level representations with the original image, and an adversarial loss that encourages the reconstructed image to lie on the natural image manifold. Their method is able

to produce much sharper reconstructions. However, it still does not tackle the problem of uncertainty in reconstruction. Given a high-level feature representation, the deconvolutional network deterministically generates a single image, despite the fact that there exist many images having the same representation. Also, there is no obvious way to sample images from their model because the distribution of the feature representations is unknown. Concurrent to our work, [42] incorporates the feature prior with a variant of denoising auto-encoder (DAE). Their sampling relies on an iterative optimization procedure, while we are focused on efficient feed-forward sampling.

### 3. Methods

In this section we introduce our model architecture. In Section 3.1 we briefly overview the framework of Generative Adversarial Networks. We then describe our proposal for Stacked Generative Adversarial Networks in Section 3.2. In Sections 3.3 and 3.4 we will focus on our two novel loss functions, conditional loss and entropy loss, respectively.

#### 3.1. Background: Generative Adversarial Network

As shown in Figure 1 (a), the original GAN [17] is trained using a two-player min-max game: a discriminator  $D$  is trained to distinguish generated images from real images, and a generator  $G$  tries to fool the discriminator  $D$ . The discriminator loss  $\mathcal{L}_D$  and the generator loss  $\mathcal{L}_G$  are formally defined as follows:

$$\mathcal{L}_D = \mathbb{E}_{x \sim P_{data}}[-\log D(x)] + \mathbb{E}_{z \sim P_z}[-\log(1 - D(G(z)))] \quad (1)$$

$$\mathcal{L}_G = \mathbb{E}_{z \sim P_z}[-\log(D(G(z)))] \quad (2)$$

In practice,  $D$  and  $G$  are usually updated alternately. The training process implicitly matches the generated image distribution  $P_G(x)$  with the real image distribution  $P_{data}(x)$  in the training set. In other words, The adversarial training forces  $G$  to generate images that reside on the manifold of natural images.

In the original GAN framework, the single noise vector  $z$  has to encapsulate every detail about an image. As a result,  $z$  cannot focus on representing high-level invariant features because invariant features are by definition not sufficient to generate an image with enough details. Intuitively, the total variations of images could be decomposed into multiple levels, with higher-level semantic variations (*e.g.*, attributes, object categories, rough shapes) and lower-level spatial variations (*e.g.*, detailed contours and textures, background clutters). A better approach would be to use different levels of noise variables to represent different levels of variations. Although [54, 70] have tried to feed multi-scale noise into different layers of the generator, there is

nothing to teach their models to use noise variables at different levels to represent different levels of variations.

#### 3.2. Stacked Generative Adversarial Networks

**Pre-trained Encoder.** We first consider a bottom-up DNN pre-trained for classification, which is referred to as the encoder  $E$  throughout. We define a stack of bottom-up deterministic nonlinear mappings:  $h_{i+1} = E_i(h_i)$ , where  $i \in \{0, 1, \dots, N-1\}$ ,  $E_i$  consists of a sequence of neural layers (*e.g.*, convolution, pooling),  $N$  is the number of hierarchies (stacks),  $h_i (i \neq 0, N)$  are intermediate representations,  $h_N = y$  is the classification result, and  $h_0 = x$  is the input image. Note that in our formulation, each  $E_i$  can contain multiple layers and the way of grouping layers together into  $E_i$  is determined by us. The number of stacks  $N$  is therefore less than the number of layers in  $E$  and is also determined by us.

**Stacked Generators.** Provided with a pre-trained encoder  $E$ , our goal is to train a top-down generator  $G$  that inverts  $E$ . Specifically,  $G$  consists of a top-down *stack* of generators  $G_i$ , each trained to invert a bottom-up mapping  $E_i$ . Each  $G_i$  takes in a noise prior  $z_i$  as input, and produces the corresponding representation  $\hat{h}_i$ .

Since  $E_i$  is usually a many-to-one mapping, there can be many  $h_i$ s such that  $h_{i+1} = E_i(h_i)$ . To aid the uncertainty of generator  $G_i$ , we provide each  $G_i$  with the conditional information of the higher-level representation  $h_{i+1}$ , defined by the pre-trained encoder  $E$ . In other words,  $\hat{h}_i = G_i(h_{i+1}, z_i)$ .  $G_i$  therefore implicitly defines a distribution of  $p_{G_i}(\hat{h}_i | h_{i+1})$ . To sample the generated representation at stack  $i$ , we first sample  $z_i$  from a simple distribution such as Gaussian, and then compute  $\hat{h}_i$  accordingly.

**Training Generator.** Each  $G_i$  is trained using the GAN approach. The training workflow is shown in Figure 1 (b). Each generator  $G_i$  is trained with a linear combination of three loss functions: adversarial loss, conditional loss, and entropy loss.

$$\mathcal{L}_{G_i} = \lambda_1 \mathcal{L}_{G_i}^{adv} + \lambda_2 \mathcal{L}_{G_i}^{cond} + \lambda_3 \mathcal{L}_{G_i}^{ent}, \quad (3)$$

where  $\mathcal{L}_{G_i}^{adv}$ ,  $\mathcal{L}_{G_i}^{cond}$ ,  $\mathcal{L}_{G_i}^{ent}$  denote adversarial loss, conditional loss, and entropy loss respectively.  $\lambda_1, \lambda_2, \lambda_3$  are the weights associated with different loss terms. In practice, we find it sufficient to set the weights such that the magnitude of different loss terms are of similar scales. No extensive hyper-parameter tuning is needed. In this section we will first introduce the adversarial loss  $\mathcal{L}_{G_i}^{adv}$ . We will then introduce  $\mathcal{L}_{G_i}^{cond}$  and  $\mathcal{L}_{G_i}^{ent}$  in Sections 3.3 and 3.4 respectively.

For each generator  $G_i$ , we introduce a *representation discriminator*  $D_i$  that distinguishes generated representations  $\hat{h}_i$ , from “real” representations  $h_i$ . Specifically, the discriminator  $D_i$  is trained with the loss function

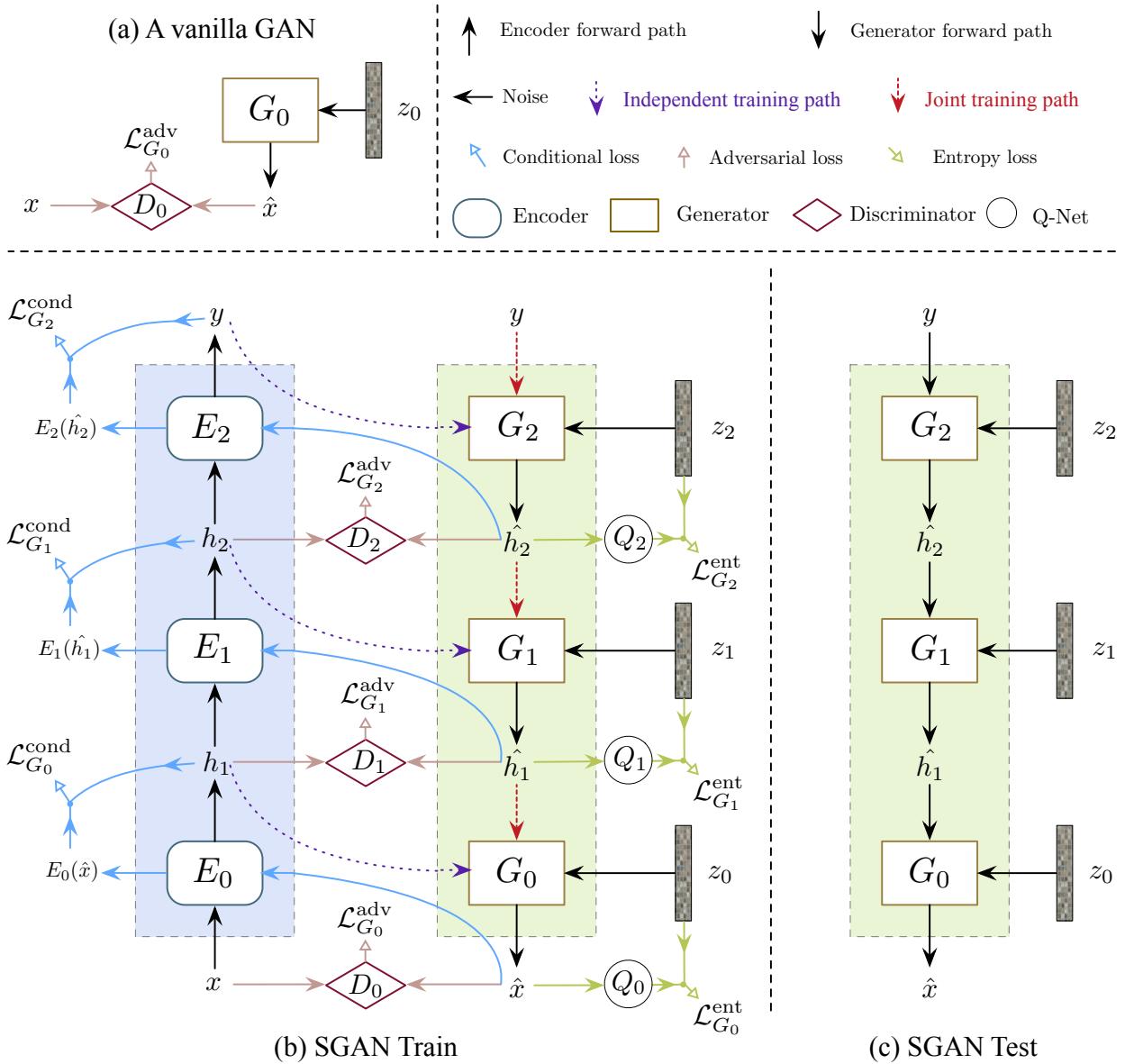


Figure 1: An overview of SGAN. (a) The original GAN in [17]. (b) The workflow of training SGAN, where each generator  $G_i$  tries to generate plausible features that can fool the corresponding representation discriminator  $D_i$ . Each generator receives conditional input from encoders in the independent training stage, and from the upper generators in the joint training stage. (c) New images can be sampled from SGAN (during test time) by feeding random noise to each generator  $G_i$ .

$$\begin{aligned} \mathcal{L}_{D_i} &= \mathbb{E}_{h_i \sim P_{data,E}}[-\log D_i(h_i)] + \\ &\mathbb{E}_{z_i \sim P_{z_i}, h_{i+1} \sim P_{data,E}}[-\log (1 - D_i(G_i(h_{i+1}, z_i)))] \quad (4) \end{aligned}$$

And  $G_i$  is trained to “fool” the representation discriminator  $D_i$ , with the adversarial loss defined by

$$\mathcal{L}_{G_i}^{adv} = \mathbb{E}_{h_{i+1} \sim P_{data,E}, z_i \sim P_{z_i}}[-\log(D_i(G_i(h_{i+1}, z_i)))] \quad (5)$$

We first train each GAN independently and then train them jointly in an end-to-end manner, as shown in Figure 1. This is similar to the training process of S<sup>2</sup>-GAN [64]. Each generator receives conditional input from encoders in the independent training stage, and from the upper generators

in the joint training stage. During joint training, the adversarial loss provided by representational discriminators can also be viewed as a type of deep supervision [35], providing intermediate supervision signals to the generator. In our current formulation,  $E$  is a discriminative model, and  $G$  is a generative model conditioned on labels. However, it is also possible to train SGAN without using label information:  $E$  can be trained with an unsupervised objective and  $G$  can be cast into an unconditional generative model by removing the label input from the top generator. We leave this for future exploration.

The idea of training a discriminator in some representation space has been explored in the context of domain adaptation [12, 2, 24], and regularization of auto-encoders [38]. However, our representation discriminators are designed to transfer the knowledge in the pre-trained bottom-up discriminative model to our top-down generative model.

**Sampling Images.** To sample images, all  $G_i$ s are stacked together in a top-down manner, as shown in Figure 1 (c). Our SGAN is capable of modeling the data distribution conditioned on the class label:  $p_G(\hat{x}|y) = p_G(\hat{h}_0|\hat{h}_N) \propto p_G(\hat{h}_0, \hat{h}_1, \dots, \hat{h}_{N-1}|\hat{h}_N) = \prod_{0 \leq i \leq N-1} p_{G_i}(\hat{h}_i|\hat{h}_{i+1})$ , where

each  $p_{G_i}(\hat{h}_i|\hat{h}_{i+1})$  is modeled by a generator  $G_i$ . From an information-theoretic perspective, SGAN factorizes the total entropy of the image distribution  $H(x)$  into multiple (smaller) conditional entropy terms:  $H(x) = H(h_0, h_1, \dots, h_N) = \sum_{i=0}^{N-1} H(h_i|h_{i+1}) + H(y)$ , thereby decomposing one difficult task into multiple easier tasks.

### 3.3. Conditional Loss

At each stack, a generator  $G_i$  is trained to capture the distribution of lower-level representations, conditioned on higher-level representations. However, in the above formulation, the generator might choose to ignore the high-level conditional information  $h_{i+1}$ , and generate plausible lower-level representations  $\hat{h}_i$  from scratch. Some previous works [40, 15, 5, 50] tackle this problem by feeding the conditional information to both the generator and discriminator. This approach, however, might introduce unnecessary complexity to the discriminator and increase model instability [47, 55].

Here we adopt a different approach: we regularize the generator by adding a loss term  $\mathcal{L}_{G_i}^{cond}$  named *conditional loss*. We feed the generated lower-level representations  $\hat{h}_i = G_i(h_{i+1}, z_i)$  back to the encoder  $E$ , and compute the recovered higher-level representations. We then enforce the recovered representations to be similar to the conditional representations. Formally:

$$\mathcal{L}_{G_i}^{cond} = \mathbb{E}_{h_{i+1} \sim P_{data, E}, z_i \sim P_{z_i}} [f(E_i(G_i(h_{i+1}, z_i)), h_{i+1})] \quad (6)$$

where  $f$  is a distance measure. We define  $f$  to be the Euclidean distance for intermediate representations and the cross-entropy for labels. Our conditional loss  $\mathcal{L}_{G_i}^{cond}$  is similar to the “feature loss” used by [7] and the “FCN loss” used by [64].  $\mathcal{L}_{G_i}^{cond}$  ensures that the generated lower-level representations are consistent with the higher-level conditional information.

### 3.4. Entropy Loss

Simply adding the conditional loss  $\mathcal{L}_{G_i}^{cond}$  leads to another issue in our experiments: the generator  $G_i$  learns to ignore the noise  $z_i$ , and compute  $\hat{h}_i$  deterministically from  $h_{i+1}$ . We refer to this phenomenon as *conditional model collapse*, *i.e.*, the conditional generative model ignores the noise input and deterministically generates the output from the conditional information. This problem has been encountered in a wide range of conditional GAN applications, *e.g.*, synthesizing future frames conditioned on previous frames [39], generating images conditioned on label maps [25], and most related to our work, synthesizing images conditioned on feature representations [7]. All the above works tried to generate *diverse* images/videos by feeding noise to the generator, but failed because the conditional generator simply ignores the noise. To our knowledge, there is still no principled way to deal with this issue. It might be tempting to think that *minibatch discrimination* [54], which encourages sample diversity in each minibatch, could solve this problem. However, even if the generator generates  $\hat{h}_i$  deterministically from  $h_{i+1}$ , the generated samples in each minibatch are still diverse since generators are conditioned on different  $h_{i+1}$ . Thus, there is no obvious way minibatch discrimination could penalize a collapsed conditional generator.

**Variational Conditional Entropy Maximization.** To alleviate this issue, we would like to encourage the generated representation  $\hat{h}_i$  to be sufficiently diverse when conditioned on  $h_{i+1}$ , *i.e.*, the conditional entropy  $H(\hat{h}_i|h_{i+1})$  should be as high as possible. Since directly maximizing  $H(\hat{h}_i|h_{i+1})$  is intractable, we propose to maximize instead a *variational lower bound* on the conditional entropy. Specifically, we use an auxiliary distribution  $Q_i(z_i|\hat{h}_i)$  to approximate the true posterior  $P_i(z_i|\hat{h}_i)$ , and augment the training objective with a loss term named *entropy loss*:

$$\mathcal{L}_{G_i}^{ent} = \mathbb{E}_{z_i \sim P_{z_i}} [\mathbb{E}_{\hat{h}_i \sim G_i(h_{i+1}, z_i)} [-\log Q_i(z_i|\hat{h}_i)]] \quad (7)$$

Below we give a proof that minimizing  $\mathcal{L}_{G_i}^{ent}$  is equivalent to maximizing a variational lower bound for  $H(\hat{h}_i|h_{i+1})$ .

$$\begin{aligned}
H(\hat{h}_i|h_{i+1}) &= H(\hat{h}_i, z_i|h_{i+1}) - H(z_i|\hat{h}_i, h_{i+1}) \\
&\geq H(\hat{h}_i, z_i|h_{i+1}) - H(z_i|\hat{h}_i) \\
&= H(z_i|h_{i+1}) + \underbrace{H(\hat{h}_i|z_i, h_{i+1})}_{0} - H(z_i|\hat{h}_i) \\
&= H(z_i|h_{i+1}) - H(z_i|\hat{h}_i) \\
&= H(z_i) - H(z_i|\hat{h}_i) \\
&= \mathbb{E}_{\hat{h}_i \sim G_i} [\mathbb{E}_{z'_i \sim P_i(z'_i|\hat{h}_i)} [\log P_i(z'_i|\hat{h}_i)]] + H(z_i) \\
&= \mathbb{E}_{\hat{h}_i \sim G_i} [\mathbb{E}_{z'_i \sim P_i(z'_i|\hat{h}_i)} [\log Q_i(z'_i|\hat{h}_i)]] \\
&\quad + \underbrace{KLD(P_i||Q_i)}_{\geq 0} + H(z_i) \\
&\geq \mathbb{E}_{\hat{h}_i \sim G_i} [\mathbb{E}_{z'_i \sim P_i(z'_i|\hat{h}_i)} [\log Q_i(z'_i|\hat{h}_i)]] + H(z_i) \\
&= \mathbb{E}_{z'_i \sim P_{z'_i}} [\mathbb{E}_{\hat{h}_i \sim G_i(\hat{h}_i|z'_i)} [\log Q_i(z'_i|\hat{h}_i)]] + H(z_i) \\
&\triangleq -\mathcal{L}_{G_i}^{ent} + H(z_i)
\end{aligned} \tag{8}$$

In practice, we parameterize  $Q_i$  with a deep network that predicts the posterior distribution of  $z_i$  given  $\hat{h}_i$ .  $Q_i$  shares most of the parameters with  $D_i$ . We treat the posterior as a diagonal Gaussian with fixed standard deviations, and use the network  $Q_i$  to only predict the posterior mean, making  $\mathcal{L}_{G_i}^{ent}$  equivalent to the Euclidean reconstruction error. In each iteration we update both  $G_i$  and  $Q_i$  to minimize  $\mathcal{L}_{G_i}^{ent}$ .

Our method is similar to the variational mutual information maximization technique proposed by [3]. A key difference is that [3] uses the  $Q$ -network to predict only a small set of deliberately constructed “latent code”, while our  $Q_i$  tries to predict *all* the noise variables  $z_i$  in each stack. The loss used in [3] therefore maximizes the *mutual information* between the output and the latent code, while ours maximizes the *conditional entropy* of the output  $\hat{h}_i$ , provided with  $h_{i+1}$ . [6, 10] also train a separate network to map images back to noise space in the context of unsupervised feature learning. Independent of our work, [4] proposes to regularize EBGAN [70] with entropy maximization, in order to prevent the discriminator from degenerating to uniform prediction. However, our entropy loss is motivated from tackling the conditional model collapse phenomenon described above.

## 4. Experiments

In the following, we perform experiments on a variety of datasets including MNIST [34], SVHN [41], and CIFAR-10 [29]. Code and pre-trained models are available at: <https://github.com/xunhuang1995/SGAN>. Readers may refer to our code repository for more details about experimental setup, hyper-parameters, *etc.*

### Encoder.

For all datasets we use a small CNN with two convolutional layers as our encoder: conv1-pool1-conv2-pool2-fc3-fc4, where fc3 is a fully connected layer and fc4 outputs classification scores before softmax. We apply horizontal flipping on CIFAR-10. No data augmentation is used on other datasets.

### Generator.

We use generators with two stacks in our experiments. Note that our framework is generally applicable to the setting with multiple stacks, and we hypothesize that using more stacks would be helpful for large-scale and high-resolution datasets. For all datasets, our top GAN  $G_1$  generates fc3 features from some random noise  $z_1$ , conditioned on label  $y$ . And the bottom GAN  $G_0$  generates images from some noise  $z_0$ , conditioned on fc3 features generated from GAN  $G_1$ . We set the loss coefficient parameters  $\lambda_1 = \lambda_2 = 1$  and  $\lambda_3 = 10$ .<sup>1</sup>

## 4.1. Datasets

We thoroughly evaluate SGAN on three widely adopted datasets: MNIST [34], SVHN [41], and CIFAR-10 [29]. The details of each dataset is described in the following.

**MNIST** The MNIST dataset contains 60,000 labeled images of hand-written digits. Each image is sized by  $28 \times 28$ .

**SVHN** The Street View House Numbers (SVHN) dataset is composed of real-world color images of house numbers collected by Google Street View [41]. Each image is of size  $32 \times 32$  and the task is to classify the digit at the center of the image. The dataset contains 73,257 images in the training set and 26,032 images in the test set.

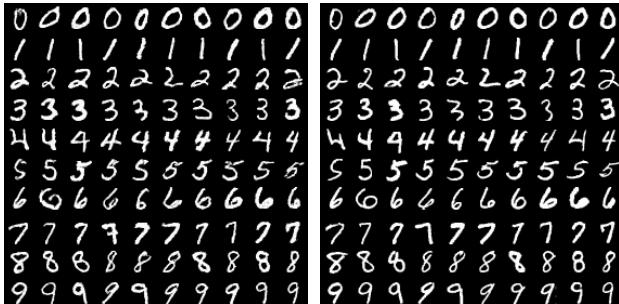
**CIFAR-10** The CIFAR-10 dataset consists of colored natural scene images sized at  $32 \times 32$  pixels. There are 50,000 training images and 10,000 test images in 10 classes.

## 4.2. Samples

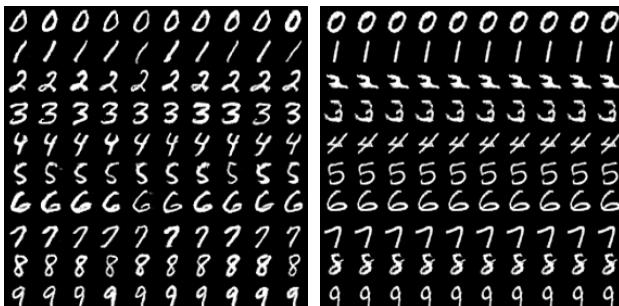
In Figure 2 (a), we show MNIST samples generated by SGAN. Each row corresponds to samples conditioned on a given digit class label. SGAN is able to generate diverse images with inherently different characteristics. The samples are visually indistinguishable from real MNIST images shown in Figure 2 (b), but still have differences compared with corresponding nearest neighbor training images.

We further examine the effect of entropy loss (Section 3.4). We show in Figure 2 (c) samples generated by bottom GAN when conditioned on a fixed fc3 feature generated by the top GAN. The samples (per row) have sufficient low-level variations, which reassures that bottom GAN learns to generate images without ignoring the noise  $z_0$ . This indicates that each GAN has the potential to capture the full conditional distribution  $p(\hat{h}_i|h_{i+1})$ . In contrast,

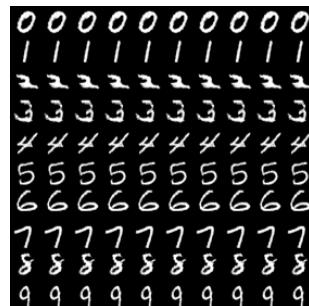
<sup>1</sup> The choice of the parameters are made so that the magnitude of each loss term is of the same scale.



(a) SGAN samples (conditioned on labels)



(c) SGAN samples (conditioned on generated fc3 features)



(d) SGAN samples (conditioned on generated fc3 features, trained without entropy loss)

**Figure 2: MNIST results.** (a) Samples generated by SGAN when conditioned on class labels. Each row corresponds to a digit class. (b) Corresponding nearest neighbor images in the MNIST training set. (c) Each row corresponds to samples generated by the bottom GAN when conditioned on a fixed  $fc_3$  feature activation, generated from the top GAN. (d) Same setting as (c), but the bottom GAN is trained without entropy loss.

in Figure 2 (d), we show samples generated without using entropy loss for bottom generator, where we observe the conditional model collapse phenomenon: the bottom GAN ignores the noise and instead deterministically generates images from  $fc_3$  features.

An advantage of SGAN compared with a vanilla GAN is its interpretability: it decomposes the total variations of an image into different levels. For example, in MNIST it decomposes the variations into  $y$  that represents the high-level digit label,  $z_1$  that captures the mid-level coarse pose of the digit and  $z_0$  that represents the low-level spatial details.

The samples generated on SVHN and CIFAR-10 datasets can also be seen in Figure 3 and Figure 4, respectively. Provided with the same  $fc_3$  feature, we see in each row of panel (c) that SGAN is able to generate samples with similar coarse outline but different lighting conditions and background clutters. Also, the nearest neighbor images in



(b) Real images (nearest neighbor labels)



(c) SGAN samples (conditioned on generated fc3 features)



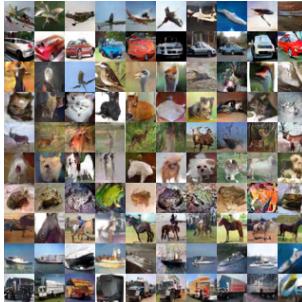
(d) SGAN samples (conditioned on generated fc3 features, trained without entropy loss)

**Figure 2: MNIST results.** (a) Samples generated by SGAN when conditioned on class labels. Each row corresponds to a digit class. (b) Corresponding nearest neighbor images in the MNIST training set. (c) Samples generated by the bottom GAN when conditioned on a fixed  $fc_3$  feature activation, generated by the top GAN. (d) Same setting as (c), but the bottom GAN is trained without entropy loss.

the training set indicate that SGAN is not simply memorizing training data, but can truly generate unseen images.

### 4.3. Comparison with the state of the art

In Table 1 we compare SGAN with other state-of-the-art generative models on CIFAR-10 dataset. The visual quality of generated images is measured by the widely used metric, Inception score [54]. Following [54], we sample 50,000 images from our model and use the code provided by [54] to compute the score. Our SGAN obtains a score of  $8.59 \pm 0.12$ , outperforming AC-GAN [44] ( $8.25 \pm 0.07$ ) and Improved GAN [54] ( $8.09 \pm 0.07$ ). Also, note that the 5 techniques introduced in [54] are not used in our implementations. Incorporating these techniques might further boost the performance of our model.



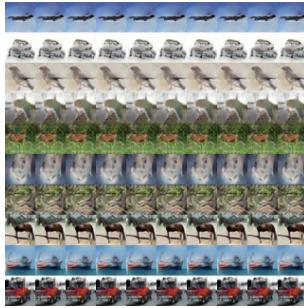
(a) SGAN samples (conditioned on labels)



(b) Real images (nearest neighbor)



(c) SGAN samples (conditioned on generated fc3 features)



(d) SGAN samples (conditioned on generated fc3 features, trained without entropy loss)

**Figure 4: CIFAR-10 results.** (a) Samples generated by SGAN when conditioned on class labels. Each row corresponds to a object category. (b) Corresponding nearest neighbor images in the CIFAR-10 training set. (c) Samples generated by the bottom GAN when conditioned on a fixed  $fc_3$  feature activation, generated by the top GAN. (d) Same setting as (c), but the bottom GAN is trained without entropy loss.

#### 4.4. More ablation studies

In Section 4.2 we have examined the effect of entropy loss on tackling the conditional model collapse problem. In order to further understand the effect of different model components, we conduct extensive ablation studies by evaluating several baseline methods on CIFAR-10 dataset. All models below use the same training hyper-parameters as the full SGAN model.

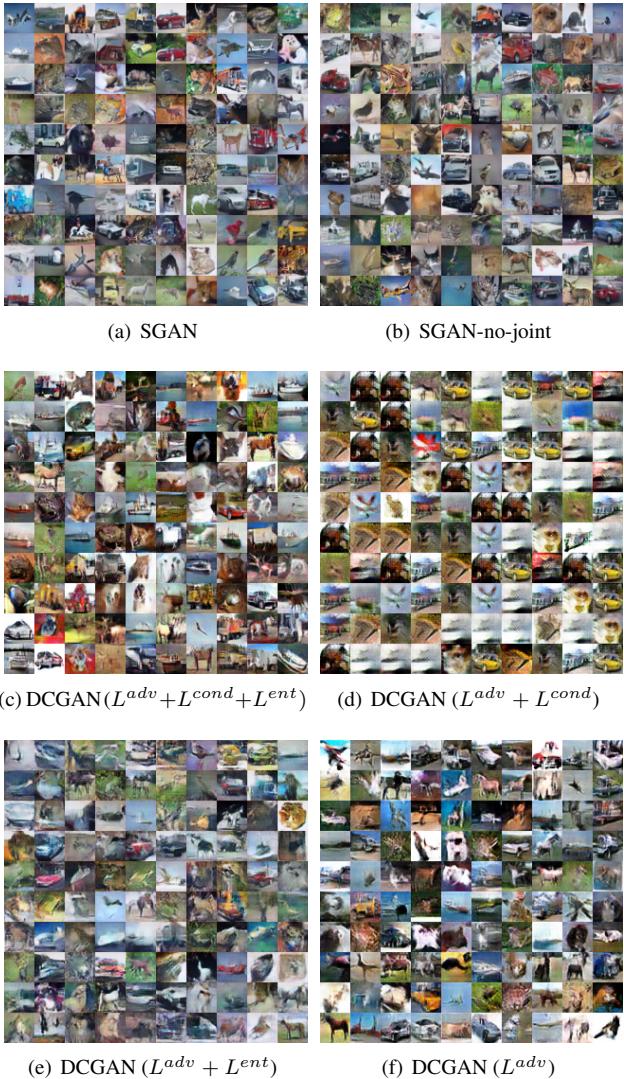
- (a) SGAN: The full model, as described in Section 3.
- (b) SGAN-no-joint: Same architecture as (a), but each GAN is trained *independently*, and there is no end-to-end joint training stage.
- (c) DCGAN ( $L^{adv} + L^{cond} + L^{ent}$ ): This is a *single* GAN model with the same architecture as the bottom GAN in SGAN, except that the generator is conditioned on

Method	Score
Infusion training [1]	$4.62 \pm 0.06$
ALI [10] (as reported in [65])	$5.34 \pm 0.05$
GMAN [11] (best variant)	$6.00 \pm 0.19$
LR-GAN [67]	$6.11 \pm 0.06$
EGAN-Ent-VI [4]	$7.07 \pm 0.10$
Denoising feature matching [65]	$7.72 \pm 0.13$
DCGAN <sup>†</sup> (with labels, as reported in [63])	$6.58$
SteinGAN <sup>†</sup> [63]	$6.35$
Improved GAN <sup>†</sup> [54] (best variant)	$8.09 \pm 0.07$
AC-GAN <sup>†</sup> [44]	$8.25 \pm 0.07$
DCGAN ( $L^{adv}$ )	$6.16 \pm 0.07$
DCGAN ( $L^{adv} + L^{ent}$ )	$5.40 \pm 0.16$
DCGAN ( $L^{adv} + L^{cond}$ ) <sup>†</sup>	$5.40 \pm 0.08$
DCGAN ( $L^{adv} + L^{cond} + L^{ent}$ ) <sup>†</sup>	$7.16 \pm 0.10$
<b>SGAN-no-joint</b> <sup>†</sup>	<b><math>8.37 \pm 0.08</math></b>
<b>SGAN</b> <sup>†</sup>	<b><math>8.59 \pm 0.12</math></b>
Real data	$11.24 \pm 0.12$

<sup>†</sup> Trained with labels.

**Table 1: Inception Score on CIFAR-10.** SGAN and SGAN-no-joint outperform state-of-the-art approaches.

- labels instead of  $fc_3$  features. Note that other techniques proposed in this paper, including conditional loss  $L^{cond}$  and entropy loss  $L^{ent}$ , are still employed. We also tried to use the full generator  $G$  in SGAN as the baseline, instead of only the bottom generator  $G_0$ . However, we failed to make it converge, possibly because  $G$  is too deep to be trained without intermediate supervision from representation discriminators.
- (d) DCGAN ( $L^{adv} + L^{cond}$ ): Same architecture as (c), but trained without entropy loss  $L^{ent}$ .
  - (e) DCGAN ( $L^{adv} + L^{ent}$ ): Same architecture as (c), but trained without conditional loss  $L^{cond}$ . The model therefore does not use label information.
  - (f) DCGAN ( $L^{adv}$ ): Same architecture as (c), but trained with neither conditional loss  $L^{cond}$  nor entropy loss  $L^{ent}$ . The model also does not use label information. It can be viewed as a plain unconditional DC-GAN model [48] and serves as the ultimate baseline.
- We show the samples and Inception scores [54] of the baseline methods in Figure 5 and Table 1, respectively. Below we summarize some of our findings:
- 1) SGAN obtains slightly higher Inception score than SGAN-no-joint. However SGAN-no-joint also generates very high quality samples and outperforms all



**Figure 5: Ablation studies on CIFAR-10.** Samples from (a) full SGAN model (b) SGAN model without joint training step. (c) DCGAN model trained with  $L^{adv} + L^{cond}$  +  $L^{ent}$  (d) DCGAN model trained with  $L^{adv} + L^{cond}$  (e) DCGAN model trained with  $L^{adv} + L^{ent}$  (f) DCGAN model trained with  $L^{adv}$ .

previous state-of-the-art models in Inception scores. We also observe that joint training is more apt to (partial) model collapse and additional caution needs to be taken, such as adding Gaussian noise to the discriminators.

- 2) SGAN, either with or without joint training, achieves significantly higher Inception score and better sample quality than the baseline DCGANs. This demonstrates the effectiveness of the proposed stacked approach.

- 3) As shown in Figure 5 (d), DCGAN ( $L^{adv} + L^{cond}$ ) collapses to generate a single image per category, while adding the entropy loss  $L^{ent}$  enables it to generate diverse images (Figure 5 (c)). This further demonstrates that entropy loss is effective at preventing conditional model collapse.
- 4) The plain unconditional DCGAN ( $L^{adv}$ ) does not collapse and obtains competitive Inception score compared with some previous models. In this case, adding entropy loss  $L^{ent}$  does not seem to provide benefits, at least in terms of the Inception score (Table 1).
- 5) The single DCGAN ( $L^{adv} + L^{cond} + L^{ent}$ ) model obtains higher Inception score than the conditional DCGAN reported in [63]. This suggests that  $L^{cond} + L^{ent}$  might offer some advantages compared to a plain conditional DCGAN, even without stacking.
- 6) In general, Inception score [54] correlates well with visual quality of images. However, it seems to be insensitive to diversity issues such as model collapse. For example, it gives the same score to Figure 5 (d) and (e) while (d) has clearly collapsed. This is consistent with results in [44, 63].

## 5. Discussion and Future Work

This paper introduces a top-down generative framework named SGAN, which effectively leverages the representational information from a pre-trained discriminative network. In contrast to GANs, our approach decomposes the hard problem of estimating image distribution into multiple relatively easier tasks – each generating plausible representations conditioned on higher-level representations. The key idea is to use representation discriminators at different training hierarchies to provide intermediate supervision. We also propose a novel entropy loss that is to our knowledge the first attempt to alleviate the conditional model collapse problem in GANs. Our entropy loss could be employed in other applications of conditional GANs, *e.g.*, synthesizing *different* future frames given the same past frames [39], or generating a *diverse* set of images conditioned on the same label map [25]. We believe this is an interesting research direction in the future.

## Acknowledgments

We would like to thank Danlu Chen for the help with Figure 1. Also, we want to thank Danlu Chen, Shuai Tang, Saining Xie, Zhuowen Tu, Felix Wu and Kilian Weinberger for helpful discussions. Yixuan Li is supported by US Army Research Office W911NF-14-1-0477. Serge Belongie is supported in part by a Google Focused Research Award.

## References

- [1] F. Bordes, S. Honari, and P. Vincent. Learning to generate samples from noise through infusion training. *ICLR submissions*, 2017. 8
- [2] X. Chen, B. Athiwaratkun, Y. Sun, K. Weinberger, and C. Cardie. Adversarial deep averaging networks for cross-lingual sentiment classification. *arXiv*, 2016. 5
- [3] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NIPS*, 2016. 6
- [4] Z. Dai, A. Almahairi, P. Bachman, E. Hovy, and A. Courville. Calibrating energy-based generative adversarial networks. *ICLR submissions*, 2017. 6, 8
- [5] E. L. Denton, S. Chintala, R. Fergus, et al. Deep generative image models using a laplacian pyramid of adversarial networks. In *NIPS*, 2015. 2, 5
- [6] J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. *ICLR submissions*, 2017. 6
- [7] A. Dosovitskiy and T. Brox. Generating images with perceptual similarity metrics based on deep networks. In *NIPS*, 2016. 2, 5
- [8] A. Dosovitskiy and T. Brox. Inverting visual representations with convolutional networks. In *CVPR*, 2016. 2
- [9] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *CVPR*, 2015. 1
- [10] V. Dumoulin, I. Belghazi, B. Poole, A. Lamb, M. Arjovsky, O. Mastropietro, and A. Courville. Adversarially learned inference. *ICLR submissions*, 2017. 6, 8
- [11] I. Durugkar, I. Gemp, and S. Mahadevan. Generative multi-adversarial networks. *ICLR submissions*, 2017. 8
- [12] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 2016. 5
- [13] L. Gatys, A. S. Ecker, and M. Bethge. Texture synthesis using convolutional neural networks. In *NIPS*, 2015. 2
- [14] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv*, 2015. 2
- [15] J. Gauthier. Conditional generative adversarial nets for convolutional face generation. *Class Project for Stanford CS231N: Convolutional Neural Networks for Visual Recognition, Winter semester*, 2014, 2014. 5
- [16] M. Germain, K. Gregor, I. Murray, and H. Larochelle. Made: masked autoencoder for distribution estimation. In *ICML*, 2015. 2
- [17] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 1, 2, 3, 4
- [18] K. Gregor, I. Danihelka, A. Graves, D. Rezende, and D. Wierstra. Draw: A recurrent neural network for image generation. In *ICML*, 2015. 2
- [19] K. Gregor, I. Danihelka, A. Mnih, C. Blundell, and D. Wierstra. Deep autoregressive networks. In *ICML*, 2014. 2
- [20] S. Gupta, J. Hoffman, and J. Malik. Cross modal distillation for supervision transfer. In *CVPR*, 2016. 1, 2
- [21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1
- [22] G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002. 2
- [23] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006. 2
- [24] J. Hoffman, D. Wang, F. Yu, and T. Darrell. Fcns in the wild: Pixel-level adversarial and constraint-based adaptation. *arxiv*, 2016. 5
- [25] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *arxiv*, 2016. 5, 9
- [26] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 2
- [27] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *NIPS*, 2014. 2
- [28] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 1, 2
- [29] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *technical report*, 2009. 6
- [30] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1
- [31] A. Lamb, V. Dumoulin, and A. Courville. Discriminative regularization for generative models. In *ICML*, 2016. 2
- [32] H. Larochelle and I. Murray. The neural autoregressive distribution estimator. In *AISTATS*, 2011. 2
- [33] A. B. L. Larsen, S. K. Sønderby, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. In *ICML*, 2016. 2
- [34] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 6
- [35] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu. Deeply-supervised nets. In *AISTATS*, 2015. 5
- [36] Y. Li, K. Swersky, and R. Zemel. Generative moment matching networks. In *ICML*, 2015. 1
- [37] A. Mahendran and A. Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *IJCV*, pages 1–23, 2016. 2
- [38] A. Makhzani, J. Shlens, N. Jaitly, and I. Goodfellow. Adversarial autoencoders. In *NIPS*, 2016. 1, 5
- [39] M. Mathieu, C. Couprie, and Y. LeCun. Deep multi-scale video prediction beyond mean square error. In *ICLR*, 2016. 5, 9
- [40] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv*, 2014. 5
- [41] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. 2011. 6
- [42] A. Nguyen, J. Yosinski, Y. Bengio, A. Dosovitskiy, and J. Clune. Plug & play generative networks: Conditional iterative generation of images in latent space. *arXiv*, 2016. 3

- [43] A. Nguyen, J. Yosinski, and J. Clune. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. *arXiv*, 2016. 2
- [44] A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier gans. *ICLR submissions*, 2017. 7, 8, 9
- [45] A. v. d. Oord, N. Kalchbrenner, and K. Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, 2016. 1, 2
- [46] A. v. d. Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu. Conditional image generation with pixelcnn decoders. In *NIPS*, 2016. 2
- [47] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 5
- [48] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016. 2, 8
- [49] A. Rasmus, M. Berglund, M. Honkala, H. Valpola, and T. Raiko. Semi-supervised learning with ladder networks. In *NIPS*, 2015. 2
- [50] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 2, 5
- [51] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014. 2
- [52] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015. 2
- [53] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 1
- [54] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *NIPS*, 2016. 2, 3, 5, 7, 8, 9
- [55] P. Sangkloy, J. Lu, C. Fang, F. Yu, and J. Hays. Scribbler: Controlling deep image synthesis with sketch and color. *arXiv*, 2016. 5
- [56] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *CVPR*, 2014. 1
- [57] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv*, 2013. 2
- [58] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 1
- [59] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 1
- [60] L. Theis and M. Bethge. Generative image modeling using spatial lstms. In *NIPS*, 2015. 2
- [61] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *ICML*, 2016. 2
- [62] H. Valpola. From neural pca to deep unsupervised learning. *Adv. in Independent Component Analysis and Learning Machines*, pages 143–171, 2015. 2
- [63] D. Wang and Q. Liu. Learning to draw samples: With application to amortized mle for generative adversarial learning. *ICLR submissions*, 2017. 8, 9
- [64] X. Wang and A. Gupta. Generative image modeling using style and structure adversarial networks. In *ECCV*, 2016. 2, 4, 5
- [65] D. Warde-Farley and Y. Bengio. Improving generative adversarial networks with denoising feature matching. *ICLR submissions*, 2017. 8
- [66] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2image: Conditional image generation from visual attributes. In *ECCV*, 2016. 2
- [67] J. Yang, A. Kannan, D. Batra, and D. Parikh. Lr-gan: Layered recursive generative adversarial networks for image generation. *ICLR submissions*, 2017. 8
- [68] Y. Zhang, K. Lee, and H. Lee. Augmenting supervised neural networks with unsupervised objectives for large-scale image classification. In *ICML*, 2016. 2
- [69] J. Zhao, M. Mathieu, R. Goroshin, and Y. Lecun. Stacked what-where auto-encoders. *ICLR Workshop*, 2016. 2
- [70] J. Zhao, M. Mathieu, and Y. LeCun. Energy-based generative adversarial network. *arXiv*, 2016. 1, 3, 6