

# Divergence Triangle for Joint Training of Generator Model, Energy-based Model, and Inference Model

Tian Han\*, Erik Nijkamp\*, Xiaolin Fang, Mitch Hill, Song-Chun Zhu, Ying Nian Wu  
Department of Statistics, UCLA

**Abstract**—This paper proposes the divergence triangle as a framework for joint training of generator model, energy-based model and inference model. The divergence triangle is a compact and symmetric (anti-symmetric) objective function that seamlessly integrates variational learning, adversarial learning, wake-sleep algorithm, and contrastive divergence in a unified probabilistic formulation. This unification makes the processes of sampling, inference, energy evaluation readily available without the need for costly Markov chain Monte Carlo methods. Our experiments demonstrate that the divergence triangle is capable of learning (1) an energy-based model with well-formed energy landscape, (2) direct sampling in the form of a generator network, and (3) feed-forward inference that faithfully reconstructs observed as well as synthesized data. The divergence triangle is a robust training method that can learn from incomplete data.

**Index Terms**—Deep generative models, Unsupervised learning, Variational inference, Adversarial contrastive divergence



## 1 INTRODUCTION

### 1.1 Integrating Three Models

Deep probabilistic generative models are a powerful framework for representing complex data distributions. They have been widely used in unsupervised learning problems to learn from unlabeled data. The goal of generative learning is to build rich and flexible models to fit complex, multi-modal data distributions as well as to be able to generate samples with high realism. The family of generative models may be roughly divided into two classes: The first class is the *energy-based model* (a.k.a undirected graphical model) and the second class is the latent variable model (a.k.a directed graphical model) which usually includes *generator model* for the generation and *inference model* for inference or reconstruction.

These models have their advantages and limitations. An energy-based model defines an explicit likelihood of the observed data up to a normalizing constant. However, sampling from such a model usually requires expensive Markov chain Monte Carlo (MCMC). A generator model defines direct sampling of the data. However, it does not have an explicit likelihood. The inference of the latent variables also requires MCMC sampling from the posterior distribution. The inference model defines an explicit approximation to the posterior distribution of the latent variables.

Combining the energy-based model, the generator model, and the inference model to get the best of each model is an attractive goal. On the other hand, challenges may accumulate when the models are trained together since different models need to effectively compete or cooperate together to achieve their highest performances. In this work, we propose the divergence triangle for joint training of energy-based model, generator model and inference model.

The learning of three models can then be seamlessly integrated in a principled probabilistic framework. The energy-based model is learned based on the samples supplied by the generator model. With the help of the inference model, the generator model is trained by both the observed data and the energy-based model. The inference model is learned from both the real data fitted by the generator model as well as the synthesized data generated by the generator model.

Our experiments demonstrate that the divergence triangle is capable of learning an energy-based model with a well-behaved energy landscape, a generator model with highly realistic samples, and an inference model with faithful reconstruction ability.

### 1.2 Prior Art

The divergence triangle jointly learns an energy-based model, a generator model, and an inference model. The following are previous methods for learning such models.

The maximum likelihood learning of the energy-based model requires expectation with respect to the current model, while the maximum likelihood learning of the generator model requires expectation with respect to the posterior distribution of the latent variables. Both expectations can be approximated by MCMC, such as Gibbs sampling [1], Langevin dynamics, or Hamiltonian Monte Carlo (HMC) [2]. [3], [4] used Langevin dynamics for learning the energy-based models, and [5] used Langevin dynamics for learning the generator model. In both cases, MCMC sampling introduces an inner loop in the training procedure, posing a computational expense.

An early version of the energy-based model is the FRAME (Filters, Random field, And Maximum Entropy) model [6], [7]. [8] used gradient-based method such as

\* Equal contributions.

Langevin dynamics to sample from the model. [9] called the energy-based models as descriptive models. [3], [4] generalized the model to deep variants.

For learning the energy-based model [10], to reduce the computational cost of MCMC sampling, contrastive divergence (CD) [11] initializes a finite step MCMC from the observed data. The resulting learning algorithm follows the gradient of the difference between two Kullback-Leibler divergences, thus the name contrastive divergence. In this paper, we shall use the term “contrastive divergence” in a more general sense than [11]. Persistent contrastive divergence [12] initializes MCMC sampling from the samples of the previous learning iteration.

Generalizing [13], [14] developed an introspective learning method where the energy function is discriminatively learned, and the energy-based model is both a generative model and a discriminative model.

For learning the generator model, the variational auto-encoder (VAE) [15], [16], [17] approximates the posterior distribution of the latent variables by an explicit inference model. In VAE, the inference model is learned jointly with the generator model from the observed data. A precursor of VAE is the wake-sleep algorithm [18], where the inference model is learned from the dream data generated by the generator model in the sleep phase.

The generator model can also be learned jointly with a discriminator model, as in the generative adversarial networks (GAN) [19], as well as deep convolutional GAN (DCGAN) [20], energy-based GAN (EB-GAN) [21], Wasserstein GAN (WGAN) [22]. GAN does not involve an inference model.

The generator model can also be learned jointly with an energy-based model [23], [24]. We can interpret the learning scheme as an adversarial version of contrastive divergence. While in GAN, the discriminator model eventually becomes a confused one, in the joint learning of the generator model and the energy-based model, the learned energy-based model becomes a well-defined probability distribution on the observed data. The joint learning bares some similarity to WGAN, but unlike WGAN, the joint learning involves two complementary probability distributions.

To bridge the gap between the generator model and the energy-based model, the cooperative learning method of [25] introduces finite-step MCMC sampling of the energy-based model with the MCMC initialized from the samples generated by the generator model. Such finite-step MCMC produces synthesized examples closer to the energy-based model, and the generator model can learn from how the finite-step MCMC revises its initial samples.

Adversarially learned inference (ALI) [26], [27] combines the learning of the generator model and inference model in an adversarial framework. ALI can be improved by adding conditional entropy regularization, resulting in the ALICE [28] model. The recently proposed method [29] shares the same spirit. They lack an energy-based model on observed data.

### 1.3 Our Contributions

Our proposed formulation, which we call the *divergence triangle*, re-interprets and integrates the following elements in

unsupervised generative learning: (1) maximum likelihood learning, (2) variational learning, (3) adversarial learning, (4) contrastive divergence, (5) wake-sleep algorithm. The learning is seamlessly integrated into a probabilistic framework based on KL divergence.

We conduct extensive experiments to analyze the learned models. Energy landscape mapping is used to verify that our learned energy-based model is well-behaved. Further, we evaluate the learning of a generator model via synthesis by generating samples with competitive fidelity, and evaluate the accuracy of the inference model both qualitatively and quantitatively via reconstruction. Our proposed model can also benefit in learning directly from incomplete images with various blocking patterns.

## 2 LEARNING DEEP PROBABILISTIC MODELS

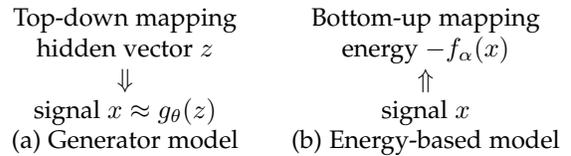
In this section, we shall review the two probabilistic models, namely the generator model and the energy-based model, both of which are parametrized by convolutional neural networks [30], [31]. Then, we shall present the maximum likelihood learning algorithms for training these two models, respectively. Our presentation of the two maximum likelihood learning algorithms is unconventional. We seek to derive both algorithms based on the Kullback-Leibler divergence using the same scheme. This will set the stage for the divergence triangle.

### 2.1 Generator Model and Energy-based Model

The generator model [15], [16], [17], [19], [20] is a generalization of the factor analysis model [32],

$$z \sim \mathcal{N}(0, I_d), \quad x = g_\theta(z) + \epsilon, \quad (1)$$

where  $g_\theta$  is a top-down mapping parametrized by a deep network with parameters  $\theta$ . It maps the  $d$ -dimensional latent vector  $z$  to the  $D$ -dimensional signal  $x$ .  $\epsilon \sim \mathcal{N}(0, \sigma^2 I_D)$  and is independent of  $z$ . In general, the model is defined by the prior distribution  $p(z)$  and the conditional distribution  $p_\theta(x|z)$ . The complete-data model  $p_\theta(z, x) = p(z)p_\theta(x|z)$ . The observed-data model is  $p_\theta(x) = \int p_\theta(z, x) dz$ . The posterior distribution is  $p_\theta(z|x) = p_\theta(z, x)/p_\theta(x)$ . See the diagram (a) below.



A complementary model is the energy-based model [3], [4], [33], [34], where  $-f_\alpha(x)$  defines the energy of  $x$ , and a low energy  $x$  is assigned a high probability. Specifically, we have the following probability model

$$\pi_\alpha(x) = \frac{1}{Z(\alpha)} \exp[f_\alpha(x)], \quad (2)$$

where  $f_\alpha(x)$  is parametrized by a bottom-up deep network with parameters  $\alpha$ , and  $Z(\alpha)$  is the normalizing constant. If  $f_\alpha(x)$  is linear in  $\alpha$ , the model becomes the familiar exponential family model in statistics or the Gibbs distribution in statistical physics. We may consider  $\pi_\alpha$  an evaluator,

where  $f_\alpha$  assigns the value to  $x$ , and  $\pi_\alpha$  evaluates  $x$  by a normalized probability distribution. See the diagram (b) above.

The energy-based model  $\pi_\alpha$  defines explicit log-likelihood via  $f_\alpha(x)$ , even though  $Z(\alpha)$  is intractable. However, it is difficult to sample from  $\pi_\alpha$ . The generator model  $p_\theta$  can generate  $x$  directly by first generating  $z \sim p(z)$ , and then transforming  $z$  to  $x$  by  $g_\theta(z)$ . But it does not define an explicit log-likelihood of  $x$ .

In the context of inverse reinforcement learning [35], [36] or inverse optimal control,  $x$  is action and  $-f_\alpha(x)$  defines the cost function or  $f_\alpha(x)$  defines the value function or the objective function.

## 2.2 Maximum Likelihood Learning

Let  $q_{\text{data}}(x)$  be the true distribution that generates the training data. Both the generator  $p_\theta$  and the energy-based model  $\pi_\alpha$  can be learned by maximum likelihood. For large sample, the maximum likelihood amounts to minimizing the Kullback-Leibler divergence  $\text{KL}(q_{\text{data}}\|p_\theta)$  over  $\theta$ , and minimizing  $\text{KL}(q_{\text{data}}\|\pi_\alpha)$  over  $\alpha$ , respectively. The expectation  $\mathbb{E}_{q_{\text{data}}}$  can be approximated by sample average.

### 2.2.1 EM-type Learning of Generator Model

To learn the generator model  $p_\theta$ , we seek to minimize  $\text{KL}(q_{\text{data}}(x)\|p_\theta(x))$  over  $\theta$ . Suppose in an iterative algorithm, the current  $\theta$  is  $\theta_t$ . We can fix  $\theta_t$  at any place we want, and vary  $\theta$  around  $\theta_t$ .

We can write

$$\begin{aligned} \text{KL}(q_{\text{data}}(x)p_{\theta_t}(z|x)\|p_\theta(z, x)) = \\ \text{KL}(q_{\text{data}}(x)\|p_\theta(x)) + \text{KL}(p_{\theta_t}(z|x)\|p_\theta(z|x)). \end{aligned} \quad (3)$$

In the EM algorithm [37], the left hand side is the surrogate objective function. This surrogate function is more tractable than the true objective function  $\text{KL}(q_{\text{data}}(x)\|p_\theta(x))$  because  $q_{\text{data}}(x)p_{\theta_t}(z|x)$  is a distribution of the complete data, and  $p_\theta(z, x)$  is the complete-data model.

We can write (3) as

$$S(\theta) = K(\theta) + \tilde{K}(\theta). \quad (4)$$

The geometric picture is that the surrogate objective function  $S(\theta)$  is above the true objective function  $K(\theta)$ , i.e.,  $S$  majorizes (upper bounds)  $K$ , and they touch each other at  $\theta_t$ , so that  $S(\theta_t) = K(\theta_t)$  and  $S'(\theta_t) = K'(\theta_t)$ . The reason is that  $\tilde{K}(\theta_t) = 0$  and  $\tilde{K}'(\theta_t) = 0$ . See Figure 1.

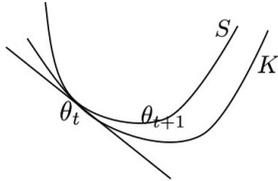


Fig. 1. The surrogate  $S$  majorizes (upper bounds)  $K$ , and they touch each other at  $\theta_t$  with the same tangent.

$q_{\text{data}}(x)p_{\theta_t}(z|x)$  gives us the complete data. Each step of EM fits the complete-data model  $p_\theta(z, x)$  by minimizing the surrogate  $S(\theta)$ ,

$$\theta_{t+1} = \arg \min_{\theta} \text{KL}(q_{\text{data}}(x)p_{\theta_t}(z|x)\|p_\theta(z, x)), \quad (5)$$

which amounts to maximizing the complete-data log-likelihood. By minimizing  $S$ , we will reduce  $S(\theta)$  relative to  $\theta_t$ , and we will reduce  $K(\theta)$  even more, relative to  $\theta_t$ , because of the majorization picture.

We can also use gradient descent to update  $\theta$ . Because  $S'(\theta_t) = K'(\theta_t)$ , and we can place  $\theta_t$  anywhere, we have

$$\begin{aligned} -\frac{\partial}{\partial \theta} \text{KL}(q_{\text{data}}(x)\|p_\theta(x)) \\ = \mathbb{E}_{q_{\text{data}}(x)p_\theta(z|x)} \left[ \frac{\partial}{\partial \theta} \log p_\theta(z, x) \right]. \end{aligned} \quad (6)$$

To implement the above updates, we need to compute the expectation with respect to the posterior distribution  $p_\theta(z|x)$ . It can be approximated by MCMC such as Langevin dynamics or HMC [2]. Both require gradient computations that can be efficiently accomplished by back-propagation. We have learned the generator using such learning method [5].

### 2.2.2 Self-critic Learning of Energy-based Model

To learn the energy-based model  $\pi_\alpha$ , we seek to minimize  $\text{KL}(q_{\text{data}}(x)\|\pi_\alpha(x))$  over  $\alpha$ . Suppose in an iterative algorithm, the current  $\alpha$  is  $\alpha_t$ . We can fix  $\alpha_t$  at any place we want, and vary  $\alpha$  around  $\alpha_t$ .

Consider the following contrastive divergence

$$\text{KL}(q_{\text{data}}(x)\|\pi_\alpha(x)) - \text{KL}(\pi_{\alpha_t}(x)\|\pi_\alpha(x)). \quad (7)$$

We can use the above as surrogate function, which is more tractable than the true objective function, since the  $\log Z(\theta)$  term is canceled out. Specifically, we can write (7) as

$$S(\alpha) = K(\alpha) - \tilde{K}(\alpha) \quad (8)$$

$$= -(\mathbb{E}_{q_{\text{data}}} [f_\alpha(x)] - \mathbb{E}_{\pi_{\alpha_t}} [f_\alpha(x)]) + \text{const.} \quad (9)$$

The geometric picture is that the surrogate function  $S(\alpha)$  is below the true objective function  $K(\alpha)$ , i.e.,  $S$  minorizes (lower bounds)  $K$ , and they touch each other at  $\alpha_t$ , so that  $S(\alpha_t) = K(\alpha_t)$ , and  $S'(\alpha_t) = K'(\alpha_t)$ . The reason is that  $\tilde{K}(\alpha_t) = 0$  and  $\tilde{K}'(\alpha_t) = 0$ . See Figure 2.

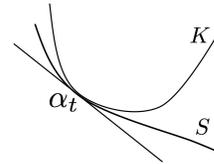


Fig. 2. The surrogate  $S$  minorizes (lower bounds)  $K$ , and they touch each other at  $\alpha_t$  with the same tangent.

Because  $S$  minorizes  $K$ , we do not have an EM-like update. However, we can still use gradient descent to update  $\alpha$ , where the derivative is

$$K'(\alpha_t) = S'(\alpha_t) = -(\mathbb{E}_{q_{\text{data}}} [f'_{\alpha_t}(x)] - \mathbb{E}_{\pi_{\alpha_t}} [f'_{\alpha_t}(x)]), \quad (10)$$

where

$$f'_{\alpha_t}(x) = \frac{\partial}{\partial \alpha} f_\alpha(x) \Big|_{\alpha_t}. \quad (11)$$

Since we can place  $\alpha_t$  anywhere, we have

$$\begin{aligned} -\frac{\partial}{\partial \alpha} \text{KL}(q_{\text{data}}(x)\|\pi_\alpha(x)) \\ = \mathbb{E}_{q_{\text{data}}} \left[ \frac{\partial}{\partial \alpha} f_\alpha(x) \right] - \mathbb{E}_{\pi_\alpha} \left[ \frac{\partial}{\partial \alpha} f_\alpha(x) \right]. \end{aligned} \quad (12)$$

To implement the above update, we need to compute the expectation with respect to the current model  $\pi_{\alpha_t}$ . It can be approximated by MCMC such as Langevin dynamics or HMC that samples from  $\pi_{\alpha_t}$ . It can be efficiently implemented by gradient computation via back-propagation. We have trained the energy-based model using such learning method [3], [4].

The above learning algorithm has an adversarial interpretation. Updating  $\alpha_t$  to  $\alpha_{t+1}$  by following the gradient of  $S(\alpha) = \text{KL}(q_{\text{data}}(x) \parallel \pi_{\alpha}(x)) - \text{KL}(\pi_{\alpha_t}(x) \parallel \pi_{\alpha}(x)) = -(\mathbb{E}_{q_{\text{data}}}[f_{\alpha}(x)] - \mathbb{E}_{\pi_{\alpha_t}}[f_{\alpha}(x)]) + \text{const}$ , we seek to decrease the first KL-divergence, while we will increase the second KL-divergence, or we seek to shift the value function  $f_{\alpha}(x)$  toward the observed data and away from the synthesized data generated from the current model. That is, the model  $\pi_{\alpha}$  criticizes its current version  $\pi_{\alpha_t}$ , i.e., the model is its own adversary or its own critic.

### 2.2.3 Similarity and Difference

In both models, at  $\theta_t$  or  $\alpha_t$ , we have  $S = K$ ,  $S' = K'$ , because  $\tilde{K} = 0$  and  $\tilde{K}' = 0$ .

The difference is that in the generator model,  $S = K + \tilde{K}$ , whereas in energy-based model,  $S = K - \tilde{K}$ .

In the generator model, if we replace the intractable  $p_{\theta_t}(z|x)$  by the inference model  $q_{\phi}(z|x)$ , we get VAE.

In energy-based model, if we replace the intractable  $\pi_{\alpha_t}(x)$  by the generator  $p_{\theta}(x)$ , we get adversarial contrastive divergence (ACD). The negative sign in front of  $\tilde{K}$  is the root of the adversarial learning.

## 3 DIVERGENCE TRIANGLE: INTEGRATING ADVERSARIAL AND VARIATIONAL LEARNING

In this section, we shall first present the divergence triangle, emphasizing its compact symmetric and anti-symmetric form. Then, we shall show that it is an re-interpretation and integration of existing methods, in particular, VAE [15], [16], [17] and ACD [23], [24].

### 3.1 Loss Function

Suppose we observe training examples  $\{x_{(i)} \sim q_{\text{data}}(x)\}_{i=1}^n$  where  $q_{\text{data}}(x)$  is the unknown data distribution.  $\pi_{\alpha}(x) \propto \exp[-f_{\alpha}(x)]$  with energy function  $-f_{\alpha}$  denotes the energy-based model with parameters  $\alpha$ . The generator model  $p(z)p_{\theta}(x|z)$  has parameters  $\theta$  and latent vector  $z$ . It is trivial to sample the latent distribution  $p(z)$  and the generative process is defined as  $z \sim p(z)$ ,  $x \sim p_{\theta}(x|z)$ .

The maximum likelihood learning algorithms for both the generator and energy-based model require MCMC sampling. We modify the maximum likelihood KL-divergences by proposing a divergence triangle criterion, so that the two models can be learned jointly without MCMC. In addition to the generator  $p_{\theta}$  and energy-based model  $\pi_{\alpha}$ , we also include an inference model  $q_{\phi}(z|x)$  in the learning scheme. Such an inference model is a key component in the variational auto-encoder [15], [16], [17]. The inference model  $q_{\phi}(z|x)$  with parameters  $\phi$  maps from the data space to latent space. In the context of EM,  $q_{\phi}(z|x)$  can be considered an imputor that imputes the missing data  $z$  to get the complete data  $(z, x)$ .

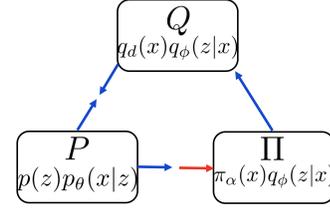


Fig. 3. Divergence triangle is based on the Kullback-Leibler divergences between three joint distributions of  $(z, x)$ . The blue arrow indicates the “running toward” behavior and the red arrow indicates the “running away” behavior.

The three models above define joint distributions over  $z$  and  $x$  from different perspectives. The two marginals, i.e., empirical data distribution  $q_{\text{data}}(x)$  and latent prior distribution  $p(z)$ , are known to us. The goal is to harmonize the three joint distributions so that the competition and cooperation between different loss terms improves learning.

The divergence triangle involves the following three joint distributions on  $(z, x)$ :

- 1)  $Q$ -distribution:  $Q(z, x) = q_{\text{data}}(x)q_{\phi}(z|x)$ .
- 2)  $P$ -distribution:  $P(z, x) = p(z)p_{\theta}(x|z)$ .
- 3)  $\Pi$ -distribution:  $\Pi(z, x) = \pi_{\alpha}(x)q_{\phi}(z|x)$ .

We propose to learn the three models  $p_{\theta}$ ,  $\pi_{\alpha}$ ,  $q_{\phi}$  by the following divergence triangle loss functional  $\mathcal{D}$

$$\max_{\alpha} \min_{\theta} \min_{\phi} \mathcal{D}(\alpha, \theta, \phi),$$

$$\mathcal{D} = \text{KL}(Q \parallel P) + \text{KL}(P \parallel \Pi) - \text{KL}(Q \parallel \Pi). \quad (13)$$

See Figure 3 for illustration. The divergence triangle is based on the three KL-divergences between the three joint distributions on  $(z, x)$ . It has a symmetric and anti-symmetric form, where the anti-symmetry is due to the negative sign in front of the last KL-divergence and the maximization over  $\alpha$ . The divergence triangle leads to the following dynamics between the three models: (1)  $Q$  and  $P$  seek to get close to each other. (2)  $P$  seeks to get close to  $\Pi$ . (3)  $\pi$  seeks to get close to  $q_{\text{data}}$ , but it seeks to get away from  $P$ , as indicated by the red arrow. Note that  $\text{KL}(Q \parallel \Pi) = \text{KL}(q_{\text{data}} \parallel \pi_{\alpha})$ , because  $q_{\phi}(z|x)$  is canceled out. The effect of (2) and (3) is that  $\pi$  gets close to  $q_{\text{data}}$ , while inducing  $P$  to get close to  $q_{\text{data}}$  as well, or in other words,  $P$  chases  $\pi_{\alpha}$  toward  $q_{\text{data}}$ .

### 3.2 Unpacking the Loss Function

The divergence triangle integrates variational and adversarial learning methods, which are modifications of maximum likelihood.

#### 3.2.1 Variational Learning

First,  $\min_{\theta} \min_{\phi} \text{KL}(Q \parallel P)$  captures the variational auto-encoder (VAE).

$$\text{KL}(Q \parallel P) = \text{KL}(q_{\text{data}}(x) \parallel p_{\theta}(x))$$

$$+ \text{KL}(q_{\phi}(z|x) \parallel p_{\theta}(z|x)), \quad (14)$$

Recall  $S = K + \tilde{K}$  in (4), if we replace the intractable  $p_{\theta_t}(z|x)$  in (4) by the explicit  $q_{\phi}(z|x)$ , we get (14), so that we avoid MCMC for sampling  $p_{\theta_t}(z|x)$ .

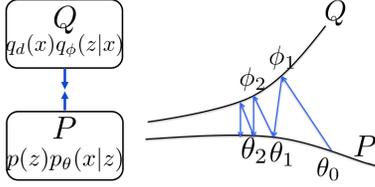


Fig. 4. Variational auto-encoder (VAE) as joint minimization by alternating projection. Left: Interaction between the models. Right: Alternating projection. The two models run toward each other.

We may interpret VAE as alternating projection between  $Q$  and  $P$ . See Figure 4 for illustration. If  $q_\phi(z|x) = p_\theta(z|x)$ , the algorithm reduces to the EM algorithm. The wake-sleep algorithm [18] is similar to VAE, except that it updates  $\phi$  by  $\min_\phi \text{KL}(P\|Q)$  instead of  $\min_\phi \text{KL}(Q\|P)$ , so that the wake-sleep algorithm does not have a single objective function.

The VAE  $\min_\theta \min_\phi \text{KL}(Q\|P)$  defines a cooperative game, with the dynamics that  $q_\phi$  and  $p_\theta$  run toward each other.

### 3.2.2 Adversarial Learning

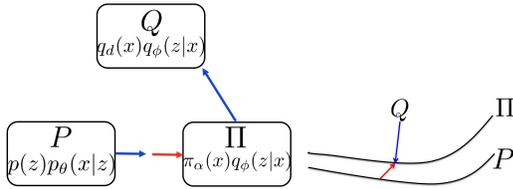


Fig. 5. Adversarial contrastive divergence (ACD). Left: Interaction between the models. Red arrow indicates a chasing game, where the generator model chases the energy-based model, which runs toward the data distribution. Right: Contrastive divergence.

Next, consider the learning of the energy-based model [23], [24]. Recall  $S = K - \tilde{K}$  in (8), if we replace the intractable  $\pi_{\alpha_t}(x)$  in (8) by  $p_\theta(x)$ , we get

$$\min_\alpha \max_\theta [\text{KL}(q_{\text{data}}(x)\|\pi_\alpha(x)) - \text{KL}(p_\theta(x)\|\pi_\alpha(x))], \quad (15)$$

or equivalently

$$\max_\alpha \min_\theta [\text{KL}(p_\theta(x)\|\pi_\alpha(x)) - \text{KL}(q_{\text{data}}(x)\|\pi_\alpha(x))], \quad (16)$$

so that we avoid MCMC for sampling  $\pi_{\alpha_t}(x)$ , and the gradient for updating  $\alpha$  becomes

$$\frac{\partial}{\partial \alpha} [\mathbb{E}_{q_{\text{data}}}(f_\alpha(x)) - \mathbb{E}_{p_\theta}(f_\alpha(x))]. \quad (17)$$

Because of the negative sign in front of the second KL-divergence in (15), we need  $\max_\theta$  in (15) or  $\min_\theta$  in (16), so that the learning becomes adversarial. See Figure 5 for illustration. Inspired by [38], we call (15) the adversarial contrastive divergence (ACD). It underlies [23], [24].

The adversarial form (15) or (16) defines a chasing game with the following dynamics: the generator  $p_\theta$  chases the energy-based model  $\pi_\alpha$  in  $\min_\theta \text{KL}(p_\theta\|\pi_\alpha)$ , the energy-based model  $\pi_\alpha$  seeks to get closer to  $q_{\text{data}}$  and get away from  $p_\theta$ . The red arrow in Figure 5 illustrates this chasing game. The result is that  $\pi_\alpha$  lures  $p_\theta$  toward  $q_{\text{data}}$ . In

the idealized case,  $p_\theta$  always catches up with  $\pi_\alpha$ , then  $\pi_\alpha$  will converge to the maximum likelihood estimate  $\min_\alpha \text{KL}(q_{\text{data}}\|\pi_\alpha)$ , and  $p_\theta$  converges to  $\pi_\alpha$ .

The above chasing game is different from VAE  $\min_\theta \min_\phi \text{KL}(Q\|P)$ , which defines a cooperative game where  $q_\phi$  and  $p_\theta$  run toward each other.

Even though the above chasing game is adversarial, both models are running toward the data distribution. While the generator model runs after the energy-based model, the energy-based model runs toward the data distribution. As a consequence, the energy-based model guides or leads the generator model toward the data distribution. It is different from GAN [19]. In GAN, the discriminator eventually becomes a confused one because the generated data become similar to the real data. In the above chasing game, the energy-based model becomes close to the data distribution.

The updating of  $\alpha$  by (17) bears similarity to Wasserstein GAN (WGAN) [22], but unlike WGAN,  $f_\alpha$  defines a probability distribution  $\pi_\alpha$ , and the learning of  $\theta$  is based on  $\min_\theta \text{KL}(p_\theta(x)\|\pi_\alpha(x))$ , which is a variational approximation to  $\pi_\alpha$ . This variational approximation only requires knowing  $f_\alpha(x)$ , without knowing  $Z(\alpha)$ . However, unlike  $q_\phi(z|x)$ ,  $p_\theta(x)$  is still intractable, in particular, its entropy does not have a closed form. Thus, we can again use variational approximation, by changing the problem to  $\min_\theta \min_\phi \text{KL}(p(z)p_\theta(x|z)\|\pi_\alpha(x)q_\phi(z|x))$ , i.e.,  $\min_\theta \min_\phi \text{KL}(P\|\Pi)$ , which is analytically tractable and which underlies [24]. In fact,

$$\text{KL}(P\|\Pi) = \text{KL}(p_\theta(x)\|\pi_\alpha(x)) + \text{KL}(p_\theta(z|x)\|q_\phi(z|x)). \quad (18)$$

Thus, we can modify (16) into  $\max_\alpha \min_\theta \min_\phi [\text{KL}(P\|\Pi) - \text{KL}(Q\|\Pi)]$ , because again  $\text{KL}(Q\|\Pi) = \text{KL}(q_{\text{data}}\|\pi_\alpha)$ .

Fitting the above together, we have the divergence triangle (13), which has a compact symmetric and anti-symmetric form.

### 3.3 Gap Between Two Models

We can write the objective function  $\mathcal{D}$  as

$$\begin{aligned} \mathcal{D} &= (\text{KL}(q_{\text{data}}(x)\|p_\theta(x)) + \text{KL}(q_\phi(z|x)\|p_\theta(z|x))) \\ &\quad - (\text{KL}(q_{\text{data}}(x)\|\pi_\alpha(x)) - \text{KL}(p(z)p_\theta(x|z)\|\pi_\alpha(x)q_\phi(z|x))) \\ &= ((\text{KL}(q_{\text{data}}(x)\|p_\theta(x)) - \text{KL}(q_{\text{data}}(x)\|\pi_\alpha(x))) \\ &\quad + \text{KL}(q_\phi(z|x)\|p_\theta(z|x)) + \text{KL}(p(z)p_\theta(x|z)\|\pi_\alpha(x)q_\phi(z|x))). \end{aligned}$$

Thus  $\mathcal{D}$  is an upper bound of the difference between the log-likelihood of the energy-based model and the log-likelihood of the generator model.

### 3.4 Two Sides of KL-divergences

In the divergence triangle, the generator model appears on the right side of  $\text{KL}(Q\|P)$ , and it also appears on the left side of  $\text{KL}(P\|\Pi)$ . The former tends to interpolate or smooth the modes of  $Q$ , while the latter tends to seek after major modes of  $\Pi$  while ignoring minor modes. As a result, the learned generator model tends to generate sharper images. As to the inference model  $q_\phi(z|x)$ , it appears on the left side of  $\text{KL}(Q\|P)$ , and it also appears on the right side of  $\text{KL}(P\|\Pi)$ . The former is variational learning of the real data, while the latter corresponds to the sleep phase of wake-sleep learning, which learns from the dream data

generated by  $P$ . The inference model thus can infer  $z$  from both observed  $x$  and generated  $x$ .

In fact, if we define

$$\mathcal{D}_0 = \text{KL}(q_{\text{data}}\|p_{\theta}) + \text{KL}(p_{\theta}\|\pi_{\alpha}) - \text{KL}(q_{\text{data}}\|\pi_{\alpha}), \quad (19)$$

we have

$$\mathcal{D} = \mathcal{D}_0 + \text{KL}(q_{\phi}(z|x)\|p_{\theta}(z|x)) + \text{KL}(p_{\theta}(z|x)\|q_{\phi}(z|x)). \quad (20)$$

(19) is the divergence triangle between the three marginal distributions on  $x$ , where  $p_{\theta}$  appears on both sides of KL-divergences. (20) is the variational scheme to make the marginal distributions into the joint distributions, which are more tractable. In (20), the two KL-divergences have reverse orders.

### 3.5 Training Algorithm

The three models are each parameterized by convolutional neural networks. The joint learning under the divergence triangle can be implemented by stochastic gradient descent, where the expectations are replaced by the sample averages. Algorithm 1 describes the procedure which is illustrated in Figure 6.

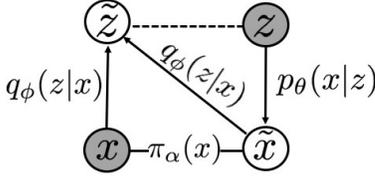


Fig. 6. Joint learning of three models. The shaded circles  $z$  and  $x$  represent variables that can be sampled from the true distributions, i.e.,  $N(0, I_d)$  and empirical data distribution, respectively.  $\tilde{x}$  and  $\tilde{z}$  are generated samples using the generator model and the inference model, respectively. The solid line with arrow represents the conditional mapping and dashed line indicates the matching loss is involved.

#### Algorithm 1 Joint Training for Divergence Triangle Model

##### Require:

- training images  $\{x_{(i)}\}_{i=1}^n$ ,
- number of learning iterations  $T$ ,
- $\alpha, \theta, \phi \leftarrow$  initialized network parameters.

##### Ensure:

- estimated parameters  $\{\alpha, \theta, \phi\}$ ,
  - generated samples  $\{\tilde{x}_{(i)}\}_{i=1}^n$ .
- 1: Let  $t \leftarrow 0$ .
  - 2: **repeat**
  - 3:  $\{z_{(i)} \sim p(z)\}_{i=1}^M$ .
  - 4:  $\{\tilde{x}_{(i)} \sim p_{\theta}(x|z_{(i)})\}_{i=1}^M$ .
  - 5:  $\{x_{(i)} \sim q_{\text{data}}(x)\}_{i=1}^M$ .
  - 6:  $\{\tilde{z}_{(i)} \sim q_{\phi}(z|x_{(i)})\}_{i=1}^M$ .
  - 7:  **$\alpha$ -step:** Given  $\{\tilde{x}_{(i)}\}_{i=1}^M$  and  $\{x_{(i)}\}_{i=1}^M$ , update  $\alpha \leftarrow \alpha + \eta_{\alpha} \frac{\partial}{\partial \alpha} \mathcal{D}$  with learning rate  $\eta_{\alpha}$ .
  - 8:  **$\phi$ -step:** Given  $\{(z_{(i)}, \tilde{x}_{(i)})\}_{i=1}^M$  and  $\{(\tilde{z}_{(i)}, x_{(i)})\}_{i=1}^M$ , update  $\phi \leftarrow \phi - \eta_{\phi} \frac{\partial}{\partial \phi} \mathcal{D}$ , with learning rate  $\eta_{\phi}$ .
  - 9:  **$\theta$ -step:** Given  $\{(z_{(i)}, \tilde{x}_{(i)})\}_{i=1}^M$  and  $\{(\tilde{z}_{(i)}, x_{(i)})\}_{i=1}^M$ , update  $\theta \leftarrow \theta - \eta_{\theta} \frac{\partial}{\partial \theta} \mathcal{D}$ , with learning rate  $\eta_{\theta}$  (optional: multiple-step update).
  - 10: Let  $t \leftarrow t + 1$ .
  - 11: **until**  $t = T$



Fig. 7. Generated samples. Left: generated samples on CIFAR-10 dataset. Right: generated samples on CelebA dataset.

## 4 EXPERIMENTS

In this section, we demonstrate not only that the divergence triangle is capable of successfully learning an energy-based model with a well-behaved energy landscape, a generator model with highly realistic samples, and an inference model with faithful reconstruction ability, but we also show competitive performance on four tasks: image generation, test image reconstruction, energy landscape mapping, and learning from incomplete images. For image generation, we consider spatial stationary texture images, temporal stationary dynamic textures, and general object categories. We also test our model on large-scale datasets and high-resolution images.

The images are resized and scaled to  $[-1, 1]$ , no further pre-processing is needed. The network parameters are initialized with zero-mean Gaussian with standard deviation 0.02 and optimized using Adam [39]. Network weights are decayed with rate 0.0005, and batch normalization [40] is used. We refer to the Appendix for the model specifications.

### 4.1 Image Generation

In this experiment, we evaluate the visual quality of generator samples from our divergence triangle model. If the generator model is well-trained, then the obtained samples should be realistic and match the visual features and contents of training images.

#### 4.1.1 Object Generation

For object categories, we test our model on two commonly-used datasets of natural images: CIFAR-10 and CelebA [41]. For CelebA face dataset, we randomly select 9,000 images for training and another 1,000 images for testing in reconstruction task. The face images are resized to  $64 \times 64$  and CIFAR-10 images remain  $32 \times 32$ . The qualitative results of generated samples for objects are shown in Figure 7. We further evaluate our model using quantitative evaluations which are based on the Inception Score (IS) [42] for CIFAR-10 and Frechet Inception Distance (FID) [43] for CelebA faces. We generate 50,000 random samples for the computation of the inception score and 10,000 random samples for the computation of the FID score. Table 1 shows the IS and FID scores of our model compared with VAE [15], DCGAN [20], WGAN [22], CoopNet [25], CEGAN [24], ALI [26], ALICE [28].

Model	VAE [15]	DCGAN [20]	WGAN [22]	CoopNet [25]	CEGAN [24]	ALI [26]	ALICE [28]	Ours
CIFAR-10 (IS)	4.08	6.16	5.76	6.55	7.07	5.93	6.02	<b>7.23</b>
CelebA (FID)	99.09	38.39	36.36	56.57	41.89	60.29	46.14	<b>31.92</b>

TABLE 1

Sample quality evaluation. Row 1: Inception scores for CIFAR-10. Row 2: FID scores for CelebA.

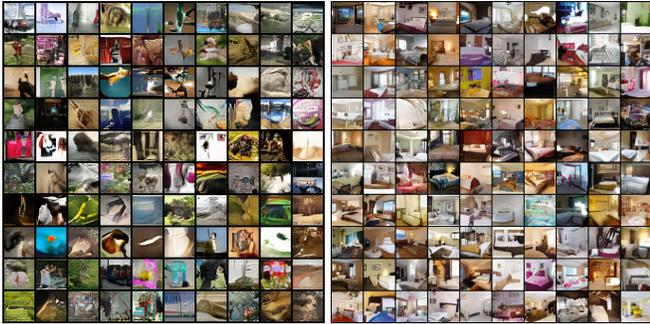


Fig. 8. Generated samples. Left:  $32 \times 32$  ImageNet. Right:  $64 \times 64$  LSUN (bedroom).

Note that for the Inception Score on CIFAR-10, we borrowed the scores from relevant papers, and for FID score on 9,000 CelebA faces, we re-implemented or used the available code with the similar network structure as our model. It can be seen that our model achieves the competitive performance compared to recent baseline models.

#### 4.1.2 Large-scale Dataset

We also train our model on large scale datasets including down-sampled  $32 \times 32$  version of ImageNet [44], [45] (roughly 1 million images) and Large-scale Scene Understand (LSUN) dataset [46]. For the LSUN dataset, we consider the *bedroom*, *tower* and *Church outdoor* categories which contains roughly 3 million, 0.7 million and 0.1 million images and were re-sized to  $64 \times 64$ . The network structures are similar with the ones used in object generation with twice the number of channels and batch normalization is used in all three models. Generated samples are shown on Figure 8.

#### 4.1.3 High-resolution Synthesis

In this section, we recruit a layer-wise training scheme to learn models on CelebA-HQ [47] with resolutions of up to  $1,024 \times 1,024$  pixels. Layer-wise training dates back to initializing deep neural networks by Restricted Boltzmann Machines to overcome optimization hurdles [48], [49] and has been resurrected in progressive GANs [47], albeit the order of layer transitions is reversed such that top layers are trained first. This resembles a Laplacian Pyramid [50] in which images are generated in a coarse-to-fine fashion.

As in [47], the training starts with down-sampled images with a spatial resolution of  $4 \times 4$  while progressively increasing the size of the images and number of layers. All three models are grown in synchrony where  $1 \times 1$  convolutions project between RGB and feature. In contrast to [47], we do not require mini-batch discrimination to increase variation of  $g_\theta(\cdot)$  nor gradient penalty to preserve 1-Lipschitz continuity of  $f_\alpha(\cdot)$ .

Figure 9 depicts high-fidelity synthesis in a resolution of  $1,024 \times 1,024$  pixels sampled from the generator model  $g_\theta(z)$  on CelebA-HQ. Figure 10 illustrates linear interpolation in latent space (i.e.,  $(1 - \alpha) \cdot z_0 + \alpha \cdot z_1$ ), which indicates diversity in the samples.

Therefore, the joint learning in the triangle formulation is not only able to train the three models with stable optimization, but it also achieves synthesis with high fidelity.

#### 4.1.4 Texture Synthesis

We consider texture images, which are spatial stationary and contain repetitive patterns. The texture images are resized to  $224 \times 224$ . Separate models are trained on each image. We start from the latent factor of size  $7 \times 7 \times 5$  and use five convolutional-transpose layers with kernel size 4 and up-sampling factor 2 for the generator network. The layers have 512, 512, 256, 128 and 3 filters, respectively, and ReLU non-linearity between each layer is used. The inference model has the inverse or “mirror” structure of generator model except that we use convolutional layers and ReLU with leak factor 0.2. The energy-based model has three convolutional layers. The first two layers have kernel size 7 with stride 2 for 100 and 70 filters respectively, and the last layer has 30 filters with kernel size 5 and stride 1.

The representative examples are shown in Figure 11. Three texture synthesis results are obtained by sampling different latent factors from prior distribution  $p(z)$ . Notice that although we only have one texture image for training, the proposed triangle divergence model can effectively utilize the repetitive patterns, thus generating realistic texture images with different configurations.

#### 4.1.5 Dynamic Texture Synthesis

Our model can also be used for dynamic patterns which exhibit stationary regularity in the temporal domain. The training video clips are selected from Dyntex database [51] and resized to  $64 \text{ pixels} \times 64 \text{ pixels} \times 32 \text{ frames}$ . Inspired by recent work [52], [53], we adopt spatial-temporal models for dynamic patterns that are stationary in the temporal domain but non-stationary in the spatial domain. Specifically, we start from 10 latent factors of size  $1 \times 1 \times 2$  for each video clip and we adopt the same spatial-temporal convolutional transpose generator network as in [53] except we use kernel size 5 for the second layer. For the inference model, we use 5 spatial-temporal convolutional layers. The first 4 layers have kernel size 4 with upsampling factor 2 and the last layer is fully-connected in spatial domain but convolutional in the temporal domain, yielding re-parametrized  $\mu_\phi$  and  $\sigma_\phi$  which have the same size the as latent factors. For the energy-based model, we use three spatial-temporal convolutional layers. The first two layers have kernel size 4 with up-sample factor 2 in all directions, but the last layer is fully-connected in the spatial domain but convolutional with

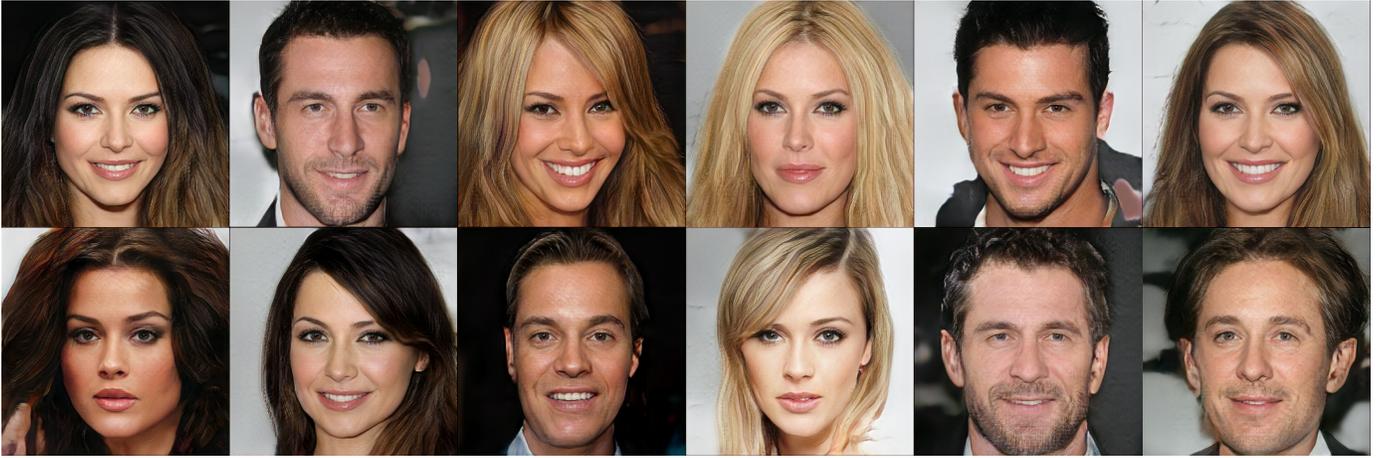


Fig. 9. Generated samples with  $1,024 \times 1,024$  resolution drawn from  $g_\theta(z)$  with 512-dimensional latent vector  $z \sim N(0, I_d)$  for CelebA-HQ.

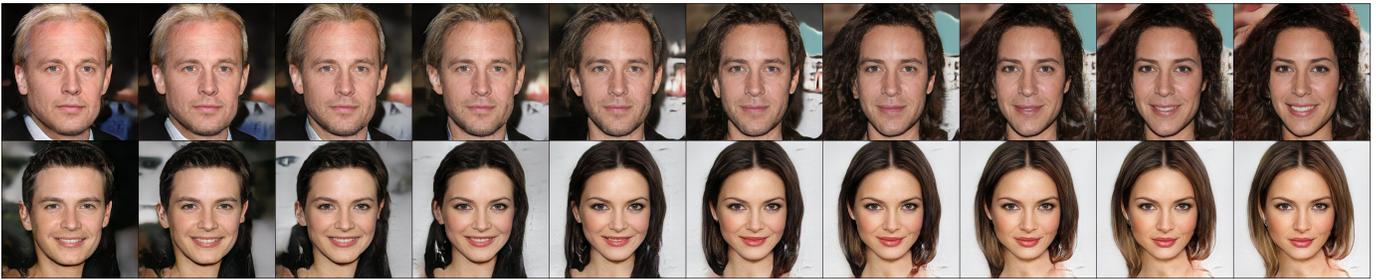


Fig. 10. High-resolution synthesis from the generator model  $g_\theta(z)$  with linear interpolation in latent space (i.e.,  $(1 - \alpha) \cdot z_0 + \alpha \cdot z_1$ ) for CelebA-HQ.



Fig. 11. Generated texture patterns. For each row, the left one is the training texture, the remaining images are 3 textures generated by divergence triangle.

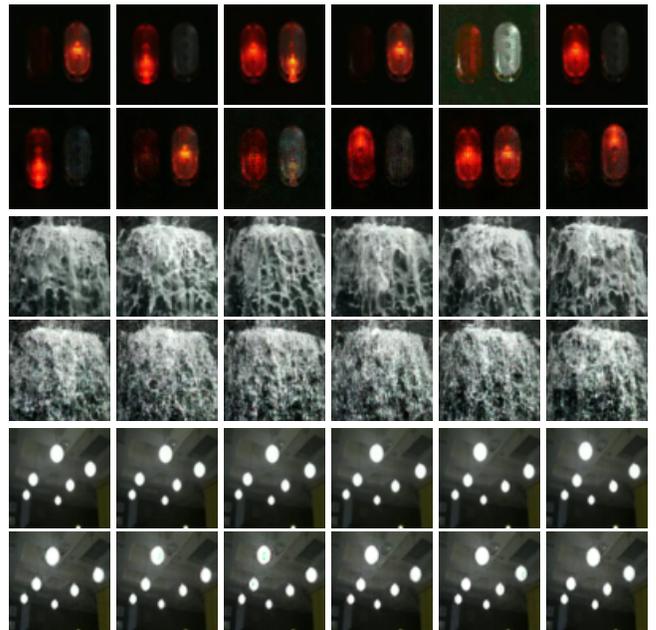


Fig. 12. Generated dynamic texture patterns. The top row shows the frames from the training video, the bottom row represents the frames for the generated video.

Model	WS [18]	VAE [15]	ALI [26]	ALICE [28]	Ours
CIFAR-10	0.058	0.037	0.311	0.034	<b>0.028</b>
CelebA	0.152	0.039	0.519	0.046	<b>0.030</b>

TABLE 2

Test reconstruction evaluation. Row 1: MSE for CIFAR-10 test set. Row 2: MSE for 1,000 hold out set from CelebA.

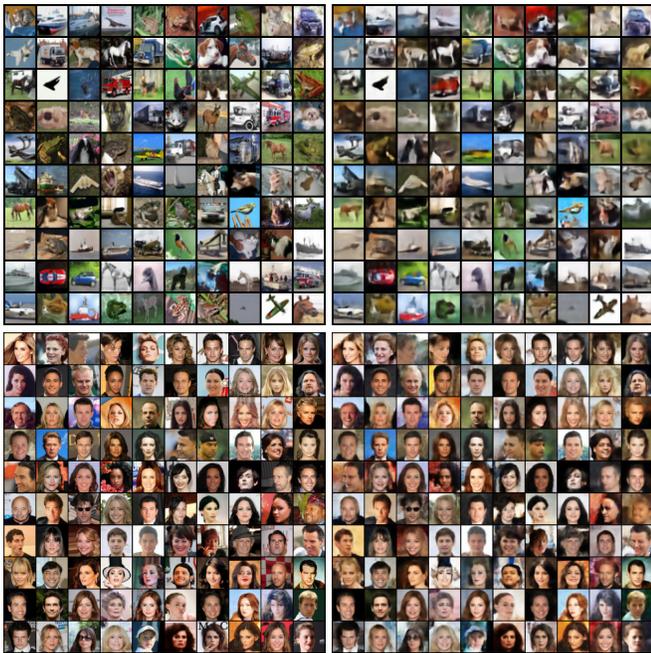


Fig. 13. Test image reconstruction. Top: CIFAR-10. Bottom: CelebA. Left: test images. Right: reconstructed images.

kernel size 4 and upsample by 2 in the temporal domain. Each layer has 64, 128 and 128 filters, respectively. Some of the synthesis results are shown in Figure 12. Note, we subsampled 6 frames of the training and generated video clips and we only show them in the first batch for illustration.

## 4.2 Test Image Reconstruction

In this experiment, we evaluate the reconstruction ability of our model for a hold-out testing image dataset. This is a strong indicator for the accuracy of our inference model. Specifically, if our divergence triangle model  $\mathcal{D}$  is well-learned, then the inference model should match the true posterior of generator model, i.e.,  $q_\phi(z|x) \approx p_\theta(z|x)$ . Therefore, given test signal  $x_{te}$ , its reconstruction  $\tilde{x}_{te}$  should be close to  $x_{te}$ , i.e.,  $x_{te} \xrightarrow{q_\phi} z_{te} \xrightarrow{p_\theta} \tilde{x}_{te} \approx x_{te}$ . Figure 13 shows the testing images and their reconstructions on CIFAR-10 and CelebA.

For CIFAR-10, we use its own 10,000 test images while for CelebA, we use the hold-out 1,000 test images as stated above. The reconstruction quality is further measured by per-pixel mean square error (MSE). Table 2 shows the per-pixel MSE of our model compared to WS [18], VAE [15], ALI [26], ALICE [28].

Note, we do not consider methods without inference models on training data, including variants of GANs and cooperative training, since it is infeasible to test such models using image reconstruction.

## 4.3 Energy Landscape Mapping

In the following, we evaluate the learned energy-based model by mapping the macroscopic structure of the energy landscape. When following a MLE regime by minimizing  $\text{KL}(q_{\text{data}} \parallel \pi_\alpha)$ , we expect the energy-function  $-f_\alpha(x)$  to encode  $x \sim q_{\text{data}}(x)$  as local energy minima. Moreover,  $-f_\alpha(x)$  should form minima for unseen images and macroscopic landscape structure in which basins of minima are distinctly separated by energy barriers. Hopfield observed that such landscape is a model of associative memory [54].

In order to learn a well-formed energy-function, in Algorithm 1, we perform multiple  $\theta$ -steps such that the samples  $\{\tilde{x}_i\}_{i=1}^M$  are sufficiently “close” to the local minima of  $-f_\alpha(x)$ . This avoids the formation of energy minima not resembling the data. The variational approximation of entropy of the marginal generator distribution  $H(p_\theta(x))$  preserves diversity in the samples avoiding mode-collapse.

To verify that (i) local minima of  $-f_\alpha(x)$  resemble  $\{x_i\}$  and (ii) minima are separated by significant energy barriers, we shall follow the approach used in [55]. When clustering with respect to energetic barriers, the landscape is partitioned into Hopfield basins of attraction whereby each point  $\{x_i\}$  on the landscape  $-f_\alpha(x)$  is mapped onto a local minimum  $\{\hat{x}_i\}$  by a steepest-descent path  $x_i^{t+1} = x_i^t + \eta \nabla f_\alpha(x_i^t)$ . The similarity measure used for hierarchical clustering is the barrier energy that separates any two regions. Given a pair of local minima  $\{\hat{x}_i, \hat{x}_j\}$ , we estimate the barrier  $b_{i,j} = \max\{-f_\alpha(x_k) : x_k \in \hat{x}_i \xrightarrow{\gamma} \hat{x}_j\}$  as the highest energy along a linear interpolation  $x \xrightarrow{\gamma} y = \{x + \gamma(y - x) : \gamma \in [0, 1]\}$ . If  $b_{i,j} < \epsilon$  for some energy threshold  $\epsilon$ , then  $\{x_i, x_j\}$  belong to the same basin. The clustering is repeated recursively until all minima are clustered together. Such graphs have come to be referred as disconnectivity graphs (DG) [56].

We conduct energy landscape mapping experiments on the MNIST [57] and Fashion-MNIST [58] datasets, each containing 70,000 grayscale images of size  $28 \times 28$  pixels depicting handwritten digits and fashion products from 10 categories, respectively. The energy landscape mapping is not without limitations, because it is practically impossible to locate all local modes. Based on the local modes located by our algorithm, see Figure 14 for the MNIST dataset, it suggests that the learned energy function is well-formed which not only encodes meaningful images as minima, but also forms meaningful macroscopic structure. Moreover, within basins the local minima have a high degree of purity (i.e. digits within a basin belong to the same class), and, the energy barrier between basins seem informative (i.e. basins of ones and sixes form pure super-basins). Figure 15 depicts the energy landscape mapping on Fashion-MNIST.

Potential applications include unsupervised classification in which energy barriers act as a geodesic similarity measure which captures perceptual distance (as opposed to e.g.  $\ell_2$  distance), weakly-supervised classification with one label per basins, or, reconstruction of incomplete data (i.e. Hopfield content-addressable memory or image inpainting).

## 4.4 Learning from incomplete images

The divergence triangle can be used to learn from occluded images. This task is challenging [5], because only parts of

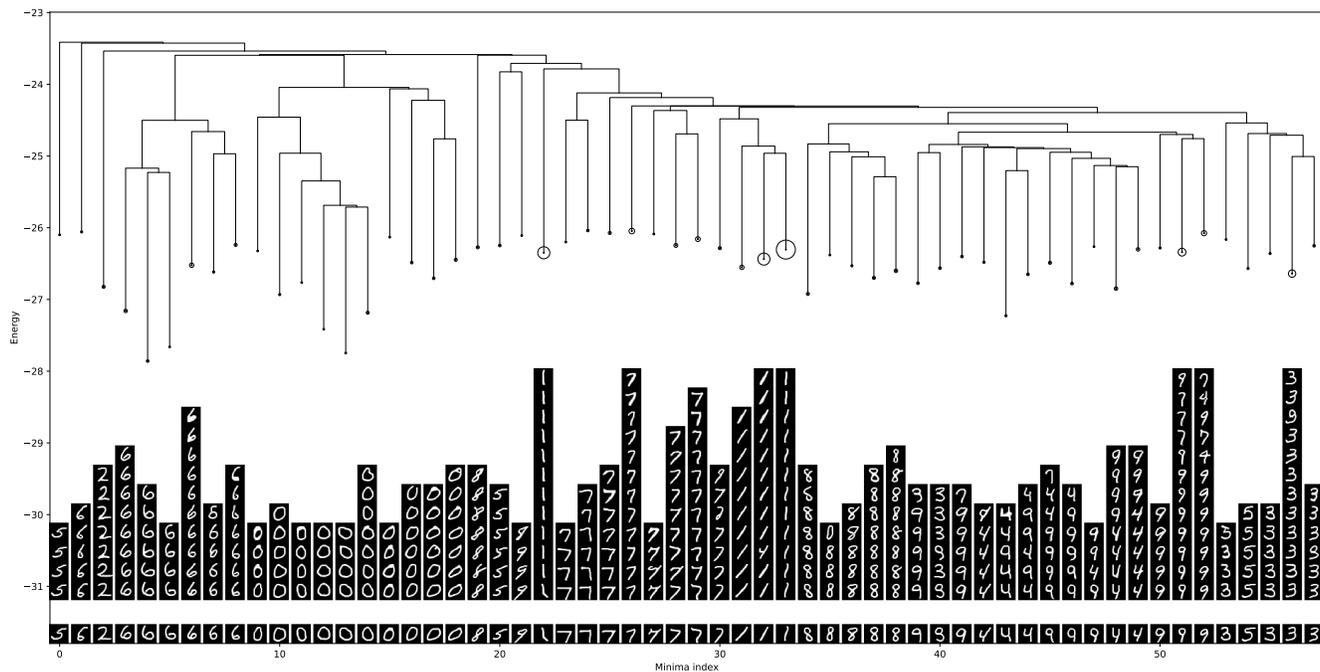


Fig. 14. Illustration of the disconnectivity-graph depicting the basin structure of the learned energy-function  $f_{\alpha}(x)$  for the MNIST dataset. Each column represents the set of at most 12 basins members ordered by energy where circles indicate the total number of basin members. Vertical lines encode minima depth in terms of energy and horizontal lines depict the lowest known barrier at which two basins merge in the landscape. Basins with less than 4 members were omitted for clarity.

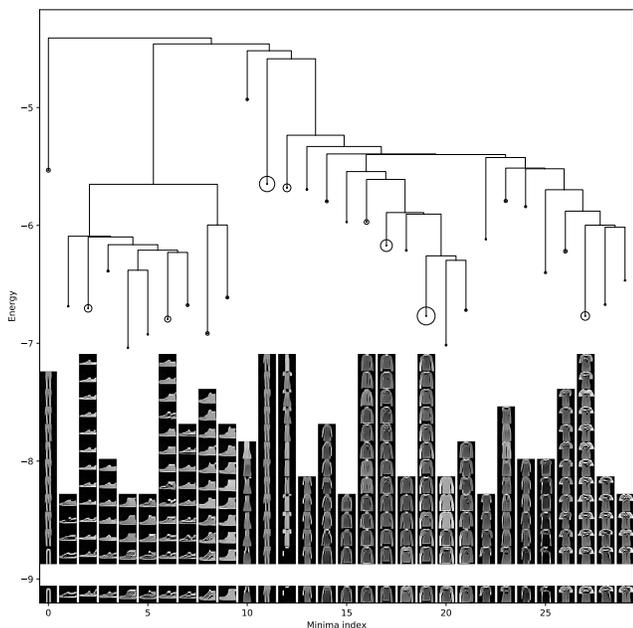


Fig. 15. Illustration of the disconnectivity-graph depicting the basin structure of the learned energy-function for the Fashion-MNIST dataset. Each column represents the set of at most 12 basins members ordered by energy where circles indicate the total number of basin members. Vertical lines encode minima depth in terms of energy and horizontal lines depict the lowest known barrier at which two basins merge in the landscape. Basins with less than 4 members were omitted for clarity.



Fig. 16. Learning from incomplete data from the CelebA dataset. The 9 columns belong to experiments P.5, P.7, MB10, MB10, B20, B20, B30, B30, B30 respectively. Row 1: original images, not observed in learning stage. Row 2: training images. Row 3: recovered images using VAE [15]. Row 4: recovered images using ABP [5]. Row 5: recovered images using our method.



Fig. 17. Image generation from different models learned from training images of the CelebA dataset with  $30 \times 30$  occlusions. Left: images generated from VAE model [15]. Middle: images generated from ABP model [5]. Right: images generated from our proposed triangle divergence model.

the images are observed, thus the model needs to learn sufficient information to recover the occluded parts. The generative models with inferential mechanism can be used for this task. Notably, [5] proposed to recover incomplete images using alternating back-propagation (ABP) which has a MCMC based inference step to refine the latent factors and perform reconstruction iteratively. VAEs [15], [59] build the inference model on occluded images, and can also be adapted for this task. It proceeds by filling the missing parts with average pixel intensity in the beginning, then iteratively re-update the missing parts using reconstructed values. Unlike VAEs, which only consider the un-occluded parts of training data, our model utilizes the generated samples which become gradually recovered during training, resulting in improved recovery accuracy and sharp generation. Note that learning from incomplete data can be difficult for variants of GANs [19], [20], [22], [24] and cooperative training [25], since inference cannot be performed directly on the occluded images.

We evaluate our model on 10,000 images randomly chosen from CelebA dataset. Then, selected images are further center cropped as in [5]. Similar to VAEs, we zero-fill the occluded parts in the beginning, then iterative update missing values using reconstructed images obtained from the generator model. Three types of occlusions are used: (1) salt and pepper noise which randomly covers 50% (P.5) and 70% (P.7) of the image. (2) Multiple block occlusion which has 10 random blocks of size  $10 \times 10$  (MB10). (3) Single block occlusion where we randomly place a large  $20 \times 20$  and  $30 \times 30$  block on each image, denoted by B20 and B30 respectively. Table 3 shows the recovery errors using VAE [15], ABP [5] and our triangle model where the error is defined as per-pixel absolute difference (relative to the range of pixel values) between the recovered image on the occluded pixels and the ground truth image.

EXP	P.5	P.7	MB10	B20	B30
VAE [15]	0.0446	0.0498	0.1169	0.0666	0.0800
ABP [5]	<b>0.0379</b>	<b>0.0428</b>	0.1070	0.0633	0.0757
Ours	0.0380	0.0430	<b>0.1060</b>	<b>0.0621</b>	<b>0.0733</b>

TABLE 3  
Recovery errors for different occlusion masks for 10,000 images selected from the CelebA dataset.

It can be seen that our model consistently out-performs the VAE model for different occlusion patterns. For structured occlusions (i.e., multiple and single blocks), the un-occluded parts contain more meaningful configurations that will improve learning of the generator through the energy-based model, which will, in turn, generate more meaningful samples to refine our inference model. This could be verified by the superior results compared to ABP [5]. While for unstructured occlusions (i.e., salt and pepper noise), ABP achieves improved recovery, a possible reason being that un-occluded parts contain less meaningful patterns which offer limited help for learning the generator and inference model. Our model synthesizes sharper and more realistic images from the generator on occluded images. See Figure 17 in which images are occluded with  $30 \times 30$  random blocks.

## 5 CONCLUSION

The proposed probabilistic framework, namely divergence triangle, for joint learning of the energy-based model, the generator model, and the inference model. The divergence triangle forms the compact learning functional for three models and naturally unifies aspects of maximum likelihood estimation [5], [25], variational auto-encoder [15], [16], [17], adversarial learning [23], [24], contrastive divergence [11], and the wake-sleep algorithm [18].

An extensive set of experiments demonstrated learning of a well-behaved energy-based model, realistic generator model as well as an accurate inference model. Moreover, experiments showed that the proposed divergence framework can be effective in learning directly from incomplete data.

In future work, we aim to extend the formulation to learn interpretable generator and energy-based models with multiple layers of sparse or semantically meaningful latent variables or features [60], [61]. Further, it would be desirable to unify the generator, energy-based and inference models into a single model [62], [63] by allowing them to share parameters and nodes instead of having separate sets of parameters and nodes.

## ACKNOWLEDGMENTS

The work is supported by DARPA XAI project N66001-17-2-4029; ARO project W911NF1810296; and ONR MURI project N00014-16-1-2007; and Extreme Science and Engineering Discovery Environment (XSEDE) grant ASC170063. We thank Dr. Tianfu Wu, Shuai Zhu and Bo Pang for helpful discussions.

## APPENDIX

### MODEL ARCHITECTURE

We describe the basic network structures, in particular for object generation. We use the following notation:

- conv( $n$ ): convolutional operation with  $n$  output feature maps.
- convT( $n$ ): convolutional transpose operation with  $n$  output feature maps.
- LReLU: Leaky-ReLU nonlinearity with default leaky factor 0.2.
- BN: Batch normalization.

The structures for CelebA (where 9,000 random images are chosen) are shown in Table 4. The structures for CIFAR-10 and MNIST/Fashion-MNIST are shown in Table 5 and Table 6, respectively.

Generator Model			
Layers	In-Out Size	Stride	BN
Input: Z	1x1x100		
4x4 convT(512), ReLU	4x4x512	1	Yes
4x4 convT(512), ReLU	8x8x512	2	Yes
4x4 convT(256), ReLU	16x16x256	2	Yes
4x4 convT(128), ReLU	32x32x128	2	Yes
4x4 convT(3), ReLU	64x64x3	2	No
Inference model			
Input: X	64x64x3		
4x4 conv(64), LReLU	32x32x64	2	Yes
4x4 conv(128), LReLU	16x16x128	2	Yes
4x4 conv(256), LReLU	8x8x256	2	Yes
4x4 conv(512), LReLU	4x4x512	2	Yes
4x4 conv(100), LReLU	$\mu, \sigma: 1x1x100$	1	No
Energy-based Model			
Input: X	64x64x3		
4x4 conv(64), LReLU	32x32x64	2	Yes
4x4 conv(128), LReLU	16x16x128	2	Yes
4x4 conv(256), LReLU	8x8x256	2	Yes
4x4 conv(256), LReLU	4x4x256	2	Yes
4x4 conv(1), LReLU	1x1x1	1	No

TABLE 4  
Network structures for CelebA (9,000).

Generator Model			
Layers	In-Out Size	Stride	BN
Input: Z	1x1x100		
4x4 convT(512), ReLU	4x4x512	1	Yes
4x4 convT(512), ReLU	8x8x512	2	Yes
4x4 convT(256), ReLU	16x16x256	2	Yes
4x4 convT(128), ReLU	32x32x128	2	Yes
3x3 convT(3), Tanh	32x32x3	1	No
Inference model			
Input: X	32x32x3		
3x3 conv(64), LReLU	32x32x64	1	No
4x4 conv(128), LReLU	16x16x128	2	No
4x4 conv(256), LReLU	8x8x256	2	No
4x4 conv(512), LReLU	4x4x512	2	No
4x4 conv(100)	$\mu, \sigma: 1x1x100$	1	No
Energy-based Model			
Input: X	32x32x3		
3x3 conv(64), LReLU	32x32x64	1	No
4x4 conv(128), LReLU	16x16x128	2	No
4x4 conv(256), LReLU	8x8x256	2	No
4x4 conv(256), LReLU	4x4x256	2	No
4x4 conv(1)	1x1x1	1	No

TABLE 5  
Network structures for CIFAR-10.

Generator Model			
Layers	In-Out Size	Stride	BN
Input: Z	1x1x100		
3x3 convT(1024), ReLU	3x3x1024	1	Yes
4x4 convT(512), ReLU	7x7x512	2	Yes
4x4 convT(256), ReLU	14x14x256	2	Yes
4x4 convT(1), Tanh	28x28x1	2	No
Inference model			
Input: X	28x28x1		
4x4 conv(128), LReLU	14x14x128	2	No
4x4 conv(256), LReLU	7x7x256	2	No
4x4 conv(512), LReLU	3x3x512	2	No
4x4 conv(100)	$\mu, \sigma: 1x1x100$	1	No
Energy-based Model			
Input: X	28x28x1		
4x4 conv(128), LReLU	14x14x128	2	No
4x4 conv(256), LReLU	7x7x256	2	No
4x4 conv(512), LReLU	3x3x512	2	No
4x4 conv(1)	1x1x1	1	No

TABLE 6  
Network structures for MNIST and Fashion-MNIST.

## REFERENCES

- [1] S. Geman and D. Geman, "Stochastic relaxation, gibbs distributions, and the bayesian restoration of images," *IEEE Transactions on pattern analysis and machine intelligence*, no. 6, pp. 721–741, 1984.
- [2] R. M. Neal, "Mcmc using hamiltonian dynamics," *Handbook of Markov Chain Monte Carlo*, vol. 2, 2011.
- [3] Y. Lu, S.-C. Zhu, and Y. N. Wu, "Learning FRAME models using CNN filters," in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [4] J. Xie, Y. Lu, S.-C. Zhu, and Y. N. Wu, "A theory of generative convnet," in *International Conference on Machine Learning*, 2016, pp. 2635–2644.
- [5] T. Han, Y. Lu, S.-C. Zhu, and Y. N. Wu, "Alternating back-propagation for generator network." in *AAAI*, vol. 3, 2017, p. 13.
- [6] S.-C. Zhu, Y. N. Wu, and D. Mumford, "Minimax entropy principle and its application to texture modeling," *Neural Computation*, vol. 9, no. 8, pp. 1627–1660, 1997.
- [7] Y. N. Wu, S. C. Zhu, and X. Liu, "Equivalence of julesz ensembles and frame models," *International Journal of Computer Vision*, vol. 38, no. 3, pp. 247–265, 2000.
- [8] S.-C. Zhu and D. Mumford, "Grade: Gibbs reaction and diffusion equations." in *International Conference on Computer Vision*, 1998, pp. 847–854.
- [9] S.-C. Zhu, "Statistical modeling and conceptualization of visual patterns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 6, pp. 691–712, 2003.
- [10] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang, "A tutorial on energy-based learning," *Predicting structured data*, vol. 1, no. 0, 2006.
- [11] G. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural Computation*, pp. 1771–1800, 2002.
- [12] T. Tieleman, "Training restricted boltzmann machines using approximations to the likelihood gradient," *ICML*, pp. 1064–1071, 2008.
- [13] Z. Tu, "Learning generative models via discriminative approaches," in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.
- [14] L. Jin, J. Lazarow, and Z. Tu, "Introspective learning for discriminative classification," in *Advances in Neural Information Processing Systems*, 2017.
- [15] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [16] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic back-propagation and approximate inference in deep generative models," in *International Conference on Machine Learning*, 2014, pp. 1278–1286.
- [17] A. Mnih and K. Gregor, "Neural variational inference and learning in belief networks," in *International Conference on Machine Learning*, 2014, pp. 1791–1799.
- [18] G. E. Hinton, P. Dayan, B. J. Frey, and R. M. Neal, "The "wake-sleep" algorithm for unsupervised neural networks," *Science*, vol. 268, no. 5214, pp. 1158–1161, 1995.
- [19] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [20] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.
- [21] J. Zhao, M. Mathieu, and Y. LeCun, "Energy-based generative adversarial network," *arXiv preprint arXiv:1609.03126*, 2016.
- [22] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International Conference on Machine Learning*, 2017, pp. 214–223.
- [23] T. Kim and Y. Bengio, "Deep directed generative models with energy-based probability estimation," *arXiv preprint arXiv:1606.03439*, 2016.
- [24] Z. Dai, A. Almahairi, P. Bachman, E. Hovy, and A. Courville, "Calibrating energy-based generative adversarial networks," *arXiv preprint arXiv:1702.01691*, 2017.
- [25] J. Xie, Y. Lu, R. Gao, S.-C. Zhu, and Y. N. Wu, "Cooperative training of descriptor and generator networks," *IEEE transactions on pattern analysis and machine intelligence (PAMI)*, 2018.
- [26] V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville, "Adversarially learned inference," *arXiv preprint arXiv:1606.00704*, 2016.

- [27] J. Donahue, P. Krähenbühl, and T. Darrell, "Adversarial feature learning," *arXiv preprint arXiv:1605.09782*, 2016.
- [28] C. Li, H. Liu, C. Chen, Y. Pu, L. Chen, R. Hénao, and L. Carin, "Alice: Towards understanding adversarial learning for joint distribution matching," in *Advances in Neural Information Processing Systems*, 2017, pp. 5495–5503.
- [29] L. Chen, S. Dai, Y. Pu, E. Zhou, C. Li, Q. Su, C. Chen, and L. Carin, "Symmetric variational autoencoder and connections to adversarial learning," in *International Conference on Artificial Intelligence and Statistics*, 2018, pp. 661–669.
- [30] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [31] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105.
- [32] D. B. Rubin and D. T. Thayer, "Em algorithms for ml factor analysis," *Psychometrika*, vol. 47, no. 1, pp. 69–76, 1982.
- [33] J. Ngiam, Z. Chen, P. W. Koh, and A. Y. Ng, "Learning deep energy models," in *International Conference on Machine Learning*, 2011, pp. 1105–1112.
- [34] J. Dai, Y. Lu, and Y.-N. Wu, "Generative modeling of convolutional neural networks," *arXiv preprint arXiv:1412.6296*, 2014.
- [35] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning." 2008.
- [36] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *Proceedings of the twenty-first international conference on Machine learning*, 2004, p. 1.
- [37] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the royal statistical society. Series B (methodological)*, pp. 1–38, 1977.
- [38] G. E. Hinton, "Training products of experts by minimizing contrastive divergence." *Neural Computation*, vol. 14, no. 8, pp. 1771–1800, 2002.
- [39] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [40] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [41] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [42] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *Advances in Neural Information Processing Systems*, 2016, pp. 2234–2242.
- [43] M. Lucic, K. Kurach, M. Michalski, S. Gelly, and O. Bousquet, "Are gans created equal? a large-scale study," *arXiv preprint arXiv:1711.10337*, 2017.
- [44] A. v. d. Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," *arXiv preprint arXiv:1601.06759*, 2016.
- [45] 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," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [46] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, and J. Xiao, "Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop," *arXiv preprint arXiv:1506.03365*, 2015.
- [47] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," *arXiv preprint arXiv:1710.10196*, 2017.
- [48] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [49] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," in *Advances in neural information processing systems*, 2007, pp. 153–160.
- [50] E. L. Denton, S. Chintala, R. Fergus *et al.*, "Deep generative image models using a laplacian pyramid of adversarial networks," in *NIPS*, 2015, pp. 1486–1494.
- [51] R. Péteri, S. Fazekas, and M. J. Huiskes, "Dyntex: A comprehensive database of dynamic textures," *Pattern Recognition Letters*, vol. 31, no. 12, pp. 1627–1632, 2010.
- [52] C. Vondrick, H. Pirsiavash, and A. Torralba, "Generating videos with scene dynamics," in *Advances In Neural Information Processing Systems*, 2016, pp. 613–621.
- [53] T. Han, Y. Lu, J. Wu, X. Xing, and Y. N. Wu, "Learning generator networks for dynamic patterns," in *WACV*, 2019.
- [54] J. J. Hopfield, "Neural networks and physical systems with emergent collective computational abilities," *Proceedings of the National Academy of Sciences*, vol. 79, no. 8, pp. 2554–2558, 1982.
- [55] M. Hill, E. Nijkamp, and S.-C. Zhu, "Building a telescope to look into high-dimensional image spaces," *arXiv preprint arXiv:1803.01043*, 2018.
- [56] D. J. Wales, M. A. Miller, and T. R. Walsh, "Archetypal energy landscapes," *Nature*, vol. 394, no. 6695, p. 758, 1998.
- [57] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [58] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [59] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic back-propagation and approximate inference in deep generative models," *arXiv preprint arXiv:1401.4082*, 2014.
- [60] R. Salakhutdinov and G. E. Hinton, "Deep boltzmann machines," in *International Conference on Artificial Intelligence and Statistics*, 2009.
- [61] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *International Conference on Machine Learning*, 2009, pp. 609–616.
- [62] L. Dinh, D. Krueger, and Y. Bengio, "Nice: Non-linear independent components estimation," *arXiv preprint arXiv:1410.8516*, 2014.
- [63] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using real nvp," *arXiv preprint arXiv:1605.08803*, 2016.