

Learning Latent Space Energy-Based Prior Model

Bo Pang^{*1} Tian Han^{*2} Erik Nijkamp^{*1} Song-Chun Zhu¹ Ying Nian Wu¹

¹University of California, Los Angeles ²Stevens Institute of Technology
 {bopang, enijkamp}@ucla.edu than6@stevens.edu {sczhu, ywu}@stat.ucla.edu

Abstract

We propose to learn energy-based model (EBM) in the latent space of a generator model, so that the EBM serves as a prior model that stands on the top-down network of the generator model. Both the latent space EBM and the top-down network can be learned jointly by maximum likelihood, which involves short-run MCMC sampling from both the prior and posterior distributions of the latent vector. Due to the low dimensionality of the latent space and the expressiveness of the top-down network, a simple EBM in latent space can capture regularities in the data effectively, and MCMC sampling in latent space is efficient and mixes well. We show that the learned model exhibits strong performances in terms of image and text generation and anomaly detection. *The one-page code can be found in supplementary materials.*

1 Introduction

In recent years, deep generative models have achieved impressive successes in image and text generation. A particularly simple and powerful model is the generator model [39, 22], which assumes that the observed example is generated by a low-dimensional latent vector via a top-down network, and the latent vector follows a non-informative prior distribution, such as uniform or isotropic Gaussian distribution. While we can learn an expressive top-down network to map the prior distribution to the data distribution, we can also learn an informative prior model in the latent space to further improve the expressive power of the whole model. This follows the philosophy of empirical Bayes where the prior model is learned from the observed data. Specifically, we assume the latent vector follows an energy-based model (EBM). We call this model the latent space energy-based prior model.

Both the latent space EBM and the top-down network can be learned jointly by maximum likelihood estimate (MLE). Each learning iteration involves Markov chain Monte Carlo (MCMC) sampling of the latent vector from both the prior and posterior distributions. Parameters of the prior model can then be updated based on the statistical difference between samples from the two distributions. Parameters of the top-down network can be updated based on the samples from the posterior distribution as well as the observed data.

Due to the low-dimensionality of the latent space, the energy function can be parametrized by a small multi-layer perceptron, yet the energy function can capture regularities in the data effectively because the EBM stands on an expressive top-down network. Moreover, MCMC in the latent space for both prior and posterior sampling is efficient and mixes well. Specifically, we employ short-run MCMC [55, 54, 56, 26] which runs a fixed number of steps from a fixed initial distribution. We formulate the resulting learning algorithm as a perturbation of MLE learning in terms of both objective function and estimating equation, so that the learning algorithm has a solid theoretical foundation. Within our theoretical framework, the short-run MCMC for posterior and prior sampling can also be amortized by jointly learned inference and synthesis networks. However, in this initial paper, we prefer keeping our model and learning method pure and self-contained, without mixing in learning

^{*}Equal contribution.

tricks from variational auto-encoder (VAE) [39, 62] and generative adversarial networks (GAN) [22, 60]. Thus we shall rely on short-run MCMC for simplicity. *The one-page code can be found in supplementary materials.* See our follow-up development on amortized inference in the context of semi-supervised learning [59].

We test the proposed modeling, learning and computing method on tasks such as image synthesis, text generation, as well as anomaly detection. We show that our method is competitive with prior art. See also our follow-up work on molecule generation [58].

Contributions. (1) We propose a latent space energy-based prior model that stands on the top-down network of the generator model. (2) We develop the maximum likelihood learning algorithm that learns the EBM prior and the top-down network jointly based on MCMC sampling of the latent vector from the prior and posterior distributions. (3) We further develop an efficient modification of MLE learning based on short-run MCMC sampling. (4) We provide theoretical foundation for learning based on short-run MCMC. The theoretical formulation can also be used to amortize short-run MCMC by extra inference and synthesis networks. (5) We provide strong empirical results to illustrate the proposed method.



Figure 1: Generated images for CelebA ($128 \times 128 \times 3$).

2 Model and learning

2.1 Model

Let x be an observed example such as an image or a piece of text, and let $z \in \mathbb{R}^d$ be the latent variables. The joint distribution of (x, z) is

$$p_{\theta}(x, z) = p_{\alpha}(z)p_{\beta}(x|z), \quad (1)$$

where $p_{\alpha}(z)$ is the prior model with parameters α , $p_{\beta}(x|z)$ is the top-down generation model with parameters β , and $\theta = (\alpha, \beta)$.

The prior model $p_{\alpha}(z)$ is formulated as an energy-based model,

$$p_{\alpha}(z) = \frac{1}{Z(\alpha)} \exp(f_{\alpha}(z))p_0(z). \quad (2)$$

where $p_0(z)$ is a known reference distribution, assumed to be isotropic Gaussian in this paper. $f_{\alpha}(z)$ is the negative energy and is parameterized by a small multi-layer perceptron with parameters α . $Z(\alpha) = \int \exp(f_{\alpha}(z))p_0(z)dz = E_{p_0}[\exp(f_{\alpha}(z))]$ is the normalizing constant or partition function.

The prior model (2) can be interpreted as an energy-based correction or exponential tilting of the original prior distribution p_0 , which is the prior distribution in the generator model in VAE.

The generation model is the same as the top-down network in VAE. For image modeling, assuming $x \in \mathbb{R}^D$,

$$x = g_{\beta}(z) + \epsilon, \quad (3)$$

where $\epsilon \sim N(0, \sigma^2 I_D)$, so that $p_{\beta}(x|z) \sim N(g_{\beta}(z), \sigma^2 I_D)$. As in VAE, σ^2 takes an assumed value. For text modeling, let $x = (x^{(t)}, t = 1, \dots, T)$ where each $x^{(t)}$ is a token. Following previous text VAE model [6], we define $p_{\beta}(x|z)$ as a conditional autoregressive model,

$$p_{\beta}(x|z) = \prod_{t=1}^T p_{\beta}(x^{(t)}|x^{(1)}, \dots, x^{(t-1)}, z) \quad (4)$$

which is parameterized by a recurrent network with parameters β .

In the original generator model, the top-down network g_{β} maps the unimodal prior distribution p_0 to be close to the usually highly multi-modal data distribution. The prior model in (2) refines p_0 so that

g_β maps the prior model p_α to be closer to the data distribution. The prior model p_α does not need to be highly multi-modal because of the expressiveness of g_β .

The marginal distribution is $p_\theta(x) = \int p_\theta(x, z) dz = \int p_\alpha(z) p_\beta(x|z) dz$. The posterior distribution is $p_\theta(z|x) = p_\theta(x, z)/p_\theta(x) = p_\alpha(z) p_\beta(x|z)/p_\theta(x)$.

In the above model, we exponentially tilt $p_0(z)$. We can also exponentially tilt $p_0(x, z) = p_0(z) p_\beta(x|z)$ to $p_\theta(x, z) = \frac{1}{Z(\theta)} \exp(f_\alpha(x, z)) p_0(x, z)$. Equivalently, we may also exponentially tilt $p_0(z, \epsilon) = p_0(z) p(\epsilon)$, as the mapping from (z, ϵ) to (z, x) is a change of variable. This leads to an EBM in both the latent space and data space, which makes learning and sampling more complex. Therefore, we choose to only tilt $p_0(z)$ and leave $p_\beta(x|z)$ as a directed top-down generation model.

2.2 Maximum likelihood

Suppose we observe training examples $(x_i, i = 1, \dots, n)$. The log-likelihood function is

$$L(\theta) = \sum_{i=1}^n \log p_\theta(x_i). \quad (5)$$

The learning gradient can be calculated according to

$$\nabla_\theta \log p_\theta(x) = \mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(x, z)] = \mathbb{E}_{p_\theta(z|x)} [\nabla_\theta (\log p_\alpha(z) + \log p_\beta(x|z))]. \quad (6)$$

See Theoretical derivations in the Supplementary for a detailed derivation.

For the prior model, $\nabla_\alpha \log p_\alpha(z) = \nabla_\alpha f_\alpha(z) - \mathbb{E}_{p_\alpha(z)} [\nabla_\alpha f_\alpha(z)]$. Thus the learning gradient for an example x is

$$\delta_\alpha(x) = \nabla_\alpha \log p_\theta(x) = \mathbb{E}_{p_\theta(z|x)} [\nabla_\alpha f_\alpha(z)] - \mathbb{E}_{p_\alpha(z)} [\nabla_\alpha f_\alpha(z)]. \quad (7)$$

The above equation has an empirical Bayes nature. $p_\theta(z|x)$ is based on the empirical observation x , while p_α is the prior model. α is updated based on the difference between z inferred from empirical observation x , and z sampled from the current prior.

For the generation model,

$$\delta_\beta(x) = \nabla_\beta \log p_\theta(x) = \mathbb{E}_{p_\theta(z|x)} [\nabla_\beta \log p_\beta(x|z)], \quad (8)$$

where $\log p_\beta(x|z) = -\|x - g_\beta(z)\|^2/(2\sigma^2) + \text{const}$ or $\sum_{t=1}^T \log p_\beta(x^{(t)}|x^{(1)}, \dots, x^{(t-1)}, z)$ for image and text modeling respectively.

Expectations in (7) and (8) require MCMC sampling of the prior model $p_\alpha(z)$ and the posterior distribution $p_\theta(z|x)$. We can use Langevin dynamics [51, 88]. For a target distribution $\pi(z)$, the dynamics iterates $z_{k+1} = z_k + s \nabla_z \log \pi(z_k) + \sqrt{2s} \epsilon_k$, where k indexes the time step of the Langevin dynamics, s is a small step size, and $\epsilon_k \sim \mathcal{N}(0, I_d)$ is the Gaussian white noise. $\pi(z)$ can be either $p_\alpha(z)$ or $p_\theta(z|x)$. In either case, $\nabla_z \log \pi(z)$ can be efficiently computed by back-propagation.

2.3 Short-run MCMC

Convergence of Langevin dynamics to the target distribution requires infinite steps with infinitesimal step size, which is impractical. We thus propose to use short-run MCMC [55, 54, 56] for approximate sampling.

The short-run Langevin dynamics is always initialized from the fixed initial distribution p_0 , and only runs a fixed number of K steps, e.g., $K = 20$,

$$z_0 \sim p_0(z), \quad z_{k+1} = z_k + s \nabla_z \log \pi(z_k) + \sqrt{2s} \epsilon_k, \quad k = 1, \dots, K. \quad (9)$$

Denote the distribution of z_K to be $\tilde{\pi}(z)$. Because of fixed $p_0(z)$ and fixed K and s , the distribution $\tilde{\pi}$ is well defined. In this paper, we put \sim sign on top of the symbols to denote distributions or quantities produced by short-run MCMC, and for simplicity, we omit the dependence on K and s in notation. As shown in [9], the Kullback-Leibler divergence $D_{KL}(\tilde{\pi}||\pi)$ decreases to zero monotonically as $K \rightarrow \infty$.

Specifically, denote the distribution of z_K to be $\tilde{p}_\alpha(z)$ if the target $\pi(z) = p_\alpha(z)$, and denote the distribution of z_K to be $\tilde{p}_\theta(z|x)$ if $\pi(z) = p_\theta(z|x)$. We can then replace $p_\alpha(z)$ by $\tilde{p}_\alpha(z)$ and replace

$p_\theta(z|x)$ by $\tilde{p}_\theta(z|x)$ in equations (7) and (8), so that the learning gradients in equations (7) and (8) are modified to

$$\tilde{\delta}_\alpha(x) = \mathbb{E}_{\tilde{p}_\theta(z|x)}[\nabla_\alpha f_\alpha(z)] - \mathbb{E}_{\tilde{p}_{\alpha_t}(z)}[\nabla_\alpha f_\alpha(z)], \quad (10)$$

$$\tilde{\delta}_\beta(x) = \mathbb{E}_{\tilde{p}_\theta(z|x)}[\nabla_\beta \log p_\beta(x|z)]. \quad (11)$$

We then update α and β based on (10) and (11), where the expectations can be approximated by Monte Carlo samples. See our follow-up work [59] on persistent chains for prior sampling.

2.4 Algorithm

The learning and sampling algorithm is described in Algorithm 1.

Algorithm 1: Learning latent space EBM prior via short-run MCMC.

input : Learning iterations T , learning rate for prior model η_0 , learning rate for generation model η_1 , initial parameters $\theta_0 = (\alpha_0, \beta_0)$, observed examples $\{x_i\}_{i=1}^n$, batch size m , number of prior and posterior sampling steps $\{K_0, K_1\}$, and prior and posterior sampling step sizes $\{s_0, s_1\}$.

output : $\theta_T = (\alpha_T, \beta_T)$.

for $t = 0 : T - 1$ **do**

1. **Mini-batch:** Sample observed examples $\{x_i\}_{i=1}^m$.
 2. **Prior sampling:** For each x_i , sample $z_i^- \sim \tilde{p}_{\alpha_t}(z)$ using equation (9), where the target distribution $\pi(z) = p_{\alpha_t}(z)$, and $s = s_0, K = K_0$.
 3. **Posterior sampling:** For each x_i , sample $z_i^+ \sim \tilde{p}_{\theta_t}(z|x_i)$ using equation (9), where the target distribution $\pi(z) = p_{\theta_t}(z|x_i)$, and $s = s_1, K = K_1$.
 4. **Learning prior model:** $\alpha_{t+1} = \alpha_t + \eta_0 \frac{1}{m} \sum_{i=1}^m [\nabla_\alpha f_{\alpha_t}(z_i^+) - \nabla_\alpha f_{\alpha_t}(z_i^-)]$.
 5. **Learning generation model:** $\beta_{t+1} = \beta_t + \eta_1 \frac{1}{m} \sum_{i=1}^m \nabla_\beta \log p_{\beta_t}(x_i|z_i^+)$.
-

The posterior sampling and prior sampling correspond to the positive phase and negative phase of latent EBM [1].

2.5 Theoretical understanding

The learning algorithm based on short-run MCMC sampling in Algorithm 1 is a modification or perturbation of maximum likelihood learning, where we replace $p_\alpha(z)$ and $p_\theta(z|x)$ by $\tilde{p}_\alpha(z)$ and $\tilde{p}_\theta(z|x)$ respectively. For theoretical underpinning, we should understand this perturbation in terms of objective function and estimating equation.

In terms of objective function, define the Kullback-Leibler divergence $D_{KL}(p(x)||q(x)) = \mathbb{E}_p[\log(p(x)/q(x))]$. At iteration t , with fixed $\theta_t = (\alpha_t, \beta_t)$, consider the following computationally tractable perturbation of the log-likelihood function of θ for an observation x ,

$$\tilde{l}_\theta(x) = \log p_\theta(x) - D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_\theta(z|x)) + D_{KL}(\tilde{p}_{\alpha_t}(z)||p_\alpha(z)). \quad (12)$$

The above is a function of θ , while θ_t is fixed. Then

$$\tilde{\delta}_\alpha(x) = \nabla_\alpha \tilde{l}_\theta(x), \quad \tilde{\delta}_\beta(x) = \nabla_\beta \tilde{l}_\theta(x), \quad (13)$$

where the derivative is taken at θ_t . See Theoretical derivations in the Supplementary for details. Thus the updating rule of Algorithm 1 follows the stochastic gradient (i.e., Monte Carlo approximation of the gradient) of a perturbation of the log-likelihood. Because θ_t is fixed, we can drop the entropies of $\tilde{p}_{\theta_t}(z|x)$ and $\tilde{p}_{\alpha_t}(z)$ in the above Kullback-Leibler divergences, hence the updating rule follows the stochastic gradient of

$$Q(\theta) = L(\theta) + \sum_{i=1}^n \left[\mathbb{E}_{\tilde{p}_{\theta_t}(z_i|x_i)}[\log p_\theta(z_i|x_i)] - \mathbb{E}_{\tilde{p}_{\alpha_t}(z)}[\log p_\alpha(z)] \right], \quad (14)$$

where $L(\theta)$ is the total log-likelihood defined in equation (5), and the gradient is taken at θ_t .

In equation (12), the first D_{KL} term is related to the EM algorithm [14]. It leads to the more tractable complete-data log-likelihood. The second D_{KL} term is related to contrastive divergence [69], except that the short-run MCMC for $\tilde{p}_{\alpha_t}(z)$ is initialized from $p_0(z)$. It serves to cancel the intractable $\log Z(\alpha)$ term.

In terms of estimating equation, the stochastic gradient descent in Algorithm 1 is a Robbins-Monro stochastic approximation algorithm [63] that solves the following estimating equation:

$$\frac{1}{n} \sum_{i=1}^n \tilde{\delta}_\alpha(x_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{p}_\theta(z_i|x_i)} [\nabla_\alpha f_\alpha(z_i)] - \mathbb{E}_{\tilde{p}_\alpha(z)} [\nabla_\alpha f_\alpha(z)] = 0, \quad (15)$$

$$\frac{1}{n} \sum_{i=1}^n \tilde{\delta}_\beta(x_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{p}_\theta(z_i|x_i)} [\nabla_\beta \log p_\beta(x_i|z_i)] = 0. \quad (16)$$

The solution to the above estimating equation defines an estimator of the parameters. Algorithm 1 converges to this estimator under the usual regularity conditions of Robbins-Monro [63]. If we replace $\tilde{p}_\alpha(z)$ by $p_\alpha(z)$, and $\tilde{p}_\theta(z|x)$ by $p_\theta(z|x)$, then the above estimating equation is the maximum likelihood estimating equation.

2.6 Amortized inference and synthesis

We can amortize the short-run MCMC sampling of the prior and posterior distributions of the latent vector by jointly learning an extra inference network $q_\phi(z|x)$ and an extra synthesis network $q_\psi(z)$, together with the original model. Let us re-define $\tilde{l}_\theta(x)$ in (12) by

$$\tilde{l}_{\theta,\phi,\psi}(x) = \log p_\theta(x) - D_{KL}(q_\phi(z|x) \| p_\theta(z|x)) + D_{KL}(q_\psi(z) \| p_\alpha(z)), \quad (17)$$

where we replace $\tilde{p}_{\theta_t}(z|x)$ in (12) by $q_\phi(z|x)$ and replace $\tilde{p}_{\alpha_t}(z)$ in (12) by $q_\psi(z)$. See [28, 29] for related formulations. Define $\tilde{L}(\theta, \phi, \psi) = \frac{1}{n} \sum_{i=1}^n \tilde{l}_{\theta,\phi,\psi}(x_i)$, we can jointly learn (θ, ϕ, ψ) by $\max_{\theta,\phi,\psi} \min_{\psi} \tilde{L}(\theta, \phi, \psi)$. The objective function $\tilde{L}(\theta, \phi, \psi)$ is a perturbation of the log-likelihood $L(\theta)$ in (5), where $-D_{KL}(q_\phi(z|x) \| p_\theta(z|x))$ leads to variational learning, and the learning of the inference network $q_\phi(z|x)$ follows VAE, except that we include the EBM prior $\log p_\alpha(z)$ in training $q_\phi(z|x)$ ($\log Z(\alpha)$ can be discarded as a constant relative to ϕ). The synthesis network $q_\psi(z)$ can be taken to be a flow-based model [15, 61]. $D_{KL}(q_\psi(z) \| p_\alpha(z))$ leads to adversarial training of $q_\psi(z)$ and $p_\alpha(z)$. $q_\psi(z)$ is trained as a variational approximation to $p_\alpha(z)$ (again $\log Z(\alpha)$ can be discarded as a constant relative to ψ), while $p_\alpha(z)$ is updated based on statistical difference between samples from the approximate posterior $q_\phi(z|x)$ and samples from the approximate prior $q_\psi(z)$, i.e., $p_\alpha(z)$ is a critic of $q_\psi(z)$. See supplementary materials for a formulation based on three D_{KL} terms.

In this initial paper, we prefer keeping our model and learning method clean and simple, without involving extra networks for learned computations, and without mixing in learning tricks from VAE and GAN. See our follow-up work on joint training of amortized inference network [59]. See also [80] for a temporal difference MCMC teaching scheme for amortizing MCMC.

3 Experiments

We present a set of experiments which highlight the effectiveness of our proposed model with (1) excellent synthesis for both visual and textual data outperforming state-of-the-art baselines, (2) high expressiveness of the learned prior model for both data modalities, and (3) strong performance in anomaly detection. For image data, we include SVHN [52], CelebA [46], and CIFAR-10 [40]. For text data, we include PTB [50], Yahoo [84], and SNLI [5]. We refer to the Supplementary for details. Code to reproduce the reported results is available². Recently we extend our work to construct a symbol-vector coupling model for semi-supervised learning and learn it with amortized inference for posterior inference and persistent chains for prior sampling [59], which demonstrates promising results in multiple data domains. In another followup [58], we find the latent space EBM can learn to capture complex chemical laws automatically and implicitly, enabling valid, novel, and diverse molecule generations. Besides the results detailed below, our extended experiments also corroborate our modeling strategy of building a latent space EBM for powerful generative modeling, meaningful representation learning, and stable training.

3.1 Image modeling

We evaluate the quality of the generated and reconstructed images. If the model is well-learned, the latent space EBM $\pi_\alpha(z)$ will fit the generator posterior $p_\theta(z|x)$ which in turn renders realistic

²<https://bpucla.github.io/latent-space-ebm-prior-project/>

generated samples as well as faithful reconstructions. We compare our model with VAE [39] and SRI [55] which assume a fixed Gaussian prior distribution for the latent vector and two recent strong VAE variants, 2sVAE [10] and RAE [20], whose prior distributions are learned with posterior samples in a second stage. We also compare with multi-layer generator (i.e., 5 layers of latent vectors) model [55] which admits a powerful learned prior on the bottom layer of latent vector. We follow the protocol as in [55].

Generation. The generator network p_θ in our framework is well-learned to generate samples that are realistic and share visual similarities as the training data. The qualitative results are shown in Figure 2. We further evaluate our model quantitatively by using Fréchet Inception Distance (FID) [48] in Table 1. It can be seen that our model achieves superior generation performance compared to listed baseline models.

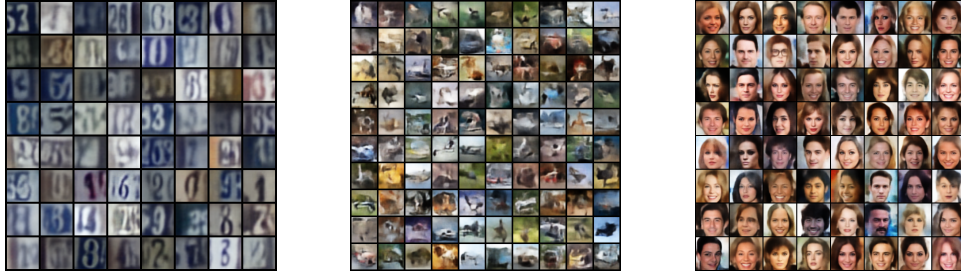


Figure 2: Generated samples for SVHN ($32 \times 32 \times 3$), CIFAR-10 ($32 \times 32 \times 3$), and CelebA ($64 \times 64 \times 3$).

Reconstruction. We evaluate the accuracy of the posterior inference by testing image reconstruction. The well-formed posterior Langevin should not only help to learn the latent space EBM model but also match the true posterior $p_\theta(z|x)$ of the generator model. We quantitatively compare reconstructions of test images with the above baseline models on mean square error (MSE). From Table 1, our proposed model could achieve not only high generation quality but also accurate reconstructions.

Models		VAE	2sVAE	RAE	SRI	SRI (L=5)	Ours
SVHN	MSE	0.019	0.019	0.014	0.018	0.011	0.008
	FID	46.78	42.81	40.02	44.86	35.23	29.44
CIFAR-10	MSE	0.057	0.056	0.027	-	-	0.020
	FID	106.37	72.90	74.16	-	-	70.15
CelebA	MSE	0.021	0.021	0.018	0.020	0.015	0.013
	FID	65.75	44.40	40.95	61.03	47.95	37.87

Table 1: MSE of testing reconstructions and FID of generated samples for SVHN ($32 \times 32 \times 3$), CIFAR-10 ($32 \times 32 \times 3$), and CelebA ($64 \times 64 \times 3$) datasets.

3.2 Text modeling

We compare our model to related baselines, SA-VAE [37], FB-VAE [45], and ARAE [86]. SA-VAE optimized posterior samples with gradient descent guided by EBLO, resembling the short run dynamics in our model. FB-VAE is the SOTA VAE for text modeling. While SA-VAE and FB-VAE assume a fixed Gaussian prior, ARAE adversarially learns a latent sample generator as an implicit prior distribution. To evaluate the quality of the generated samples, we follow [86, 8] and recruit Forward Perplexity (FPPL) and Reverse Perplexity (RPPL). FPPL is the perplexity of the generated samples evaluated under a language model trained with real data and measures the fluency of the synthesized sentences. RPPL is the perplexity of real data computed under a language model trained with the model-generated samples. Prior work employs it to measure the distributional coverage of a learned model, $p_\theta(x)$ in our case, since a model with a mode-collapsing issue results in a high RPPL. FPPL and RPPL are displayed in Table 2. Our model outperforms all the baselines on the two metrics, demonstrating the high fluency and diversity of the samples from our model. We also evaluate the reconstruction of our model against the baselines using negative log-likelihood (NLL). Our model has a similar performance as that of FB-VAE and ARAE, while they all outperform SA-VAE.

Models	FPPL	SNLI RPPL	NLL	FPPL	PTB RPPL	NLL	FPPL	Yahoo RPPL	NLL
Real Data	23.53	-	-	100.36	-	-	60.04	-	-
SA-VAE	39.03	46.43	33.56	147.92	210.02	101.28	128.19	148.57	326.70
FB-VAE	39.19	43.47	28.82	145.32	204.11	92.89	123.22	141.14	319.96
ARAE	44.30	82.20	28.14	165.23	232.93	91.31	158.37	216.77	320.09
Ours	27.81	31.96	28.90	107.45	181.54	91.35	80.91	118.08	321.18

Table 2: FPPL, RPPL, and NLL for our model and baselines on SNLI, PTB, and Yahoo datasets.

3.3 Analysis of latent space

We examine the exponential tilting of the reference prior $p_0(z)$ through Langevin samples initialized from $p_0(z)$ with target distribution $p_\alpha(z)$. As the reference distribution $p_0(z)$ is in the form of an isotropic Gaussian, we expect the energy-based correction f_α to tilt p_0 into an irregular shape. In particular, learning equation 10 may form shallow local modes for $p_\alpha(z)$. Therefore, the trajectory of a Markov chain initialized from the reference distribution $p_0(z)$ with well-learned target $p_\alpha(z)$ should depict the transition towards synthesized examples of high quality while the energy fluctuates around some constant. Figure 3 and Table 3 depict such transitions for image and textual data, respectively, which are both based on models trained with $K_0 = 40$ steps. For image data the quality of synthesis improve significantly with increasing number of steps. For textual data, there is an enhancement in semantics and syntax along the chain, which is especially clear from step 0 to 40 (see Table 3).

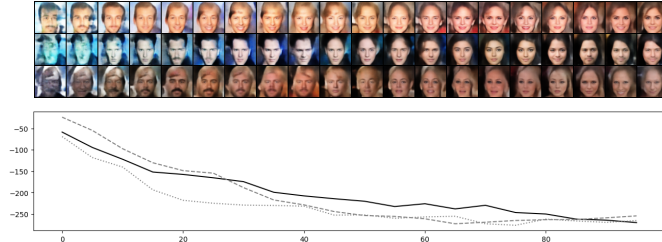


Figure 3: Transition of Markov chains initialized from $p_0(z)$ towards $\tilde{p}_\alpha(z)$ for $K'_0 = 100$ steps. *Top*: Trajectory in the CelebA data-space. *Bottom*: Energy profile over time.

judge in <unk> was not
west virginia bank <unk> which has been under N law took effect of october N
mr. peterson N years old could return to work with his clients to pay
iras must be
anticipating bonds tied to the imperial company 's revenue of \$ N million today
many of these N funds in the industrial average rose to N N from N N N
fund obtaining the the
ford 's latest move is expected to reach an agreement in principle for the sale of its loan operations
wall street has been shocked over by the merger of new york co. a world-wide financial board of the companies said it wo
n't seek strategic alternatives to the brokerage industry 's directors

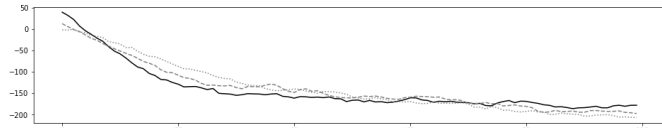


Table 3: Transition of a Markov chain initialized from $p_0(z)$ towards $\tilde{p}_\alpha(z)$. *Top*: Trajectory in the PTB data-space. Each panel contains a sample for $K'_0 \in \{0, 40, 100\}$. *Bottom*: Energy profile.

While our learning algorithm recruits short run MCMC with K_0 steps to sample from target distribution $p_\alpha(z)$, a well-learned $p_\alpha(z)$ should allow for Markov chains with realistic synthesis for $K'_0 \gg K_0$ steps. We demonstrate such long-run Markov chain with $K_0 = 40$ and $K'_0 = 2500$ in Figure 4. The long-run chain samples in the data space are reasonable and do not exhibit the oversaturating issue of the long-run chain samples of recent EBM in the data space (see oversaturating examples in Figure 3 in [54]).

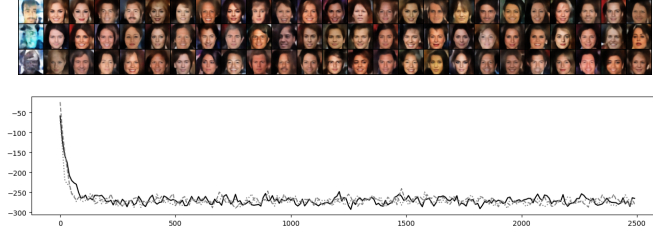


Figure 4: Transition of Markov chains initialized from $p_0(z)$ towards $\tilde{p}_\alpha(z)$ for $K'_0 = 2500$ steps. *Top*: Trajectory in the CelebA data-space for every 100 steps. *Bottom*: Energy profile over time.

3.4 Anomaly detection

We evaluate our model on anomaly detection. If the generator and EBM are well learned, then the posterior $p_\theta(z|x)$ would form a discriminative latent space that has separated probability densities for normal and anomalous data. Samples from such a latent space can then be used to detect anomalies. We take samples from the posterior of the learned model, and use the unnormalized log-posterior $\log p_\theta(x, z)$ as our decision function.

Following the protocol as in [41, 85], we make each digit class an anomaly and consider the remaining 9 digits as normal examples. Our model is trained with only normal data and tested with both normal and anomalous data. We compare with the BiGAN-based anomaly detection [85], MEG [41] and VAE using area under the precision-recall curve (AUPRC) as in [85]. Table 4 shows the results.

Heldout Digit	1	4	5	7	9
VAE	0.063	0.337	0.325	0.148	0.104
MEG	0.281 ± 0.035	0.401 ± 0.061	0.402 ± 0.062	0.290 ± 0.040	0.342 ± 0.034
BiGAN- σ	0.287 ± 0.023	0.443 ± 0.029	0.514 ± 0.029	0.347 ± 0.017	0.307 ± 0.028
Ours	0.336 ± 0.008	0.630 ± 0.017	0.619 ± 0.013	0.463 ± 0.009	0.413 ± 0.010

Table 4: AUPRC scores for unsupervised anomaly detection on MNIST. Numbers are taken from [41] and results for our model are averaged over last 10 epochs to account for variance.

3.5 Computational cost

Our method involving MCMC sampling is more costly than VAEs with amortized inference. Our model is approximately 4 times slower than VAEs on image datasets. On text datasets, ours does not have an disadvantage compared to VAEs on total training time (despite longer per-iteration time) because of better posterior samples from short run MCMC than amortized inference and the overhead of the techniques that VAEs take to address posterior collapse. To test our method’s scalability, we trained a larger generator on CelebA (128×128). It produced faithful samples (see Figure 1).

4 Discussion and conclusion

4.1 Modeling strategies and related work

We now put our work within the bigger picture of modeling and learning, and discuss related work.

Energy-based model and top-down generation model. A top-down model or a directed acyclic graphical model is of a simple factorized form that is capable of ancestral sampling. The prototype of such a model is factor analysis [65], which has been generalized to independent component analysis [35], sparse coding [57], non-negative matrix factorization [43], etc. An early example of a multi-layer top-down model is the generation model of Helmholtz machine [31]. An EBM defines an unnormalized density or a Gibbs distribution. The prototypes of such a model are exponential family distribution, the Boltzmann machine [1, 32, 66, 44], and the FRAME (Filters, Random field, And Maximum Entropy) model [89]. [87] contrasted these two classes of models, calling the top-down latent variable model the generative model, and the energy-based model the descriptive model. [24] proposed to integrate the two models, where the top-down generation model generates textons, while the EBM prior accounts for the perceptual organization or Gestalt laws of textons. Our model follows

such a plan. Recently, DVAEs [64, 76, 75] adopted restricted Boltzmann machines as the prior model for binary latent variables and a deep neural network as the top-down generation model.

The energy-based model can be translated into a classifier and vice versa via the Bayes rule [25, 73, 12, 81, 36, 42, 19, 23, 59]. The energy function in the EBM can be viewed as an objective function, a cost function, or a critic [67]. It captures regularities, rules or constraints. It is easy to specify, although optimizing or sampling the energy function requires iterative computation such as MCMC. The maximum likelihood learning of EBM can be interpreted as an adversarial scheme [83, 82, 79, 29, 17], where the MCMC serves as a generator or an actor and the energy function serves as an evaluator or a critic. The top-down generation model can be viewed as an actor [67] that directly generates samples. It is easy to sample from, though a complex top-down model is necessary for high quality samples. Comparing the two models, the scalar-valued energy function can be more expressive than the vector-valued top-down network of the same complexity, while the latter is much easier to sample from. It is thus desirable to let EBM take over the top layers of the top-down model to make it more expressive and make EBM learning feasible.

Energy-based correction of top-down model. The top-down model usually assumes independent nodes at the top layer and conditional independent nodes at subsequent layers. We can introduce energy terms at multiple layers to correct the independence or conditional independence assumptions, and to introduce inductive biases. This leads to a latent energy-based model. However, unlike undirected latent EBM, the energy-based correction is learned on top of a directed top-down model, and this can be easier than learning an undirected latent EBM from scratch. Our work is a simple example of this strategy where we correct the prior distribution. We can also correct the generation model in the data space.

From data space EBM to latent space EBM. EBM learned in data space such as image space [53, 47, 81, 18, 27, 55, 16] can be highly multi-modal, and MCMC sampling can be difficult. We can introduce latent variables and learn an EBM in latent space, while also learning a mapping from the latent space to the data space. Our work follows such a strategy. Earlier papers on this strategy are [87, 24, 4, 7, 41]. Learning EBM in latent space can be much more feasible than in data space in terms of MCMC sampling, and much of past work on EBM can be recast in the latent space.

Short-run MCMC and amortized computation. Recently, [55] proposed to use short-run MCMC to sample from the EBM in data space. [56] used it to sample the latent variables of a top-down generation model from their posterior distribution. [34] used it to improve the posterior samples from an inference network. Our work adopts short-run MCMC to sample from both the prior and the posterior of the latent variables. We provide theoretical foundation for the learning algorithm with short-run MCMC sampling. Our theoretical formulation can also be used to jointly train networks that amortize the MCMC sampling from the posterior and prior distributions.

Generator model with flexible prior. The expressive power of the generator network for image and text generation comes from the top-down network that maps a simple prior to be close to the data distribution. Most of the existing papers [49, 70, 2, 13, 74, 41] assume that the latent vector follows a given simple prior, such as isotropic Gaussian distribution or uniform distribution. However, such assumption may cause ineffective generator learning as observed in [11, 72]. Some VAE variants attempted to address the mismatch between the prior and the aggregate posterior. VampPrior [71] parameterized the prior based on the posterior inference model, while [3] proposed to construct priors using rejection sampling. ARAE [86] learned an implicit prior with adversarial training. Recently, some papers used two-stage approach [10, 20]. They first trained a VAE or deterministic auto-encoder. To enable generation from the model, they fitted a VAE or Gaussian mixture to the posterior samples from the first-stage model. VQ-VAE [77] adopted a similar approach and an autoregressive distribution over z was learned from the posterior samples. All of these prior models generally follow the empirical Bayes philosophy, which is also one motivation of our work.

4.2 Conclusion

EBM has many applications, however, its soundness and its power are limited by the difficulty with MCMC sampling. By moving from data space to latent space, and letting the EBM stand on an expressive top-down network, MCMC-based learning of EBM becomes sound and feasible, and EBM in latent space can capture regularities in data effectively. We may unleash the power of EBM in the latent space for many applications.

Broader Impact

Our work can be of interest to researchers working on generator model, energy-based models, MCMC sampling and unsupervised learning. It may also be of interest to people who are interested in image synthesis and text generation.

Acknowledgments and Disclosure of Funding

We thank the four reviewers for their insightful comments and useful suggestions. The work is supported by NSF DMS-2015577; DARPA XAI project N66001-17-2-4029; ARO project W911NF1810296; ONR MURI project N00014-16-1-2007; and XSEDE grant ASC170063. We thank the NVIDIA cooperation for the donation of 2 Titan V GPUs.

References

- [1] David H. Ackley, Geoffrey E. Hinton, and Terrence J. Sejnowski. A learning algorithm for boltzmann machines. *Cognitive Science*, 9(1):147–169, 1985.
- [2] Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 214–223, 2017.
- [3] Matthias Bauer and Andriy Mnih. Resampled priors for variational autoencoders. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 66–75, 2019.
- [4] Yoshua Bengio, Grégoire Mesnil, Yann Dauphin, and Salah Rifai. Better mixing via deep representations. In *International conference on machine learning*, pages 552–560, 2013.
- [5] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal, Sept. 2015. Association for Computational Linguistics.
- [6] Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 10–21, Berlin, Germany, Aug. 2016. Association for Computational Linguistics.
- [7] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [8] Ondřej Cífka, Aliaksei Severyn, Enrique Alfonseca, and Katja Filippova. Eval all, trust a few, do wrong to none: Comparing sentence generation models. *arXiv preprint arXiv:1804.07972*, 2018.
- [9] Thomas M. Cover and Joy A. Thomas. *Elements of information theory* (2. ed.). Wiley, 2006.
- [10] Bin Dai and David Wipf. Diagnosing and enhancing vae models. In *International Conference on Learning Representations*, 2019.
- [11] Bin Dai and David Wipf. Diagnosing and enhancing vae models. *arXiv preprint arXiv:1903.05789*, 2019.
- [12] Jifeng Dai, Yang Lu, and Ying Nian Wu. Generative modeling of convolutional neural networks. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [13] Zihang Dai, Amjad Almahairi, Philip Bachman, Eduard H. Hovy, and Aaron C. Courville. Calibrating energy-based generative adversarial networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- [14] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- [15] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- [16] Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models. *CoRR*, abs/1903.08689, 2019.
- [17] Chelsea Finn, Paul F. Christiano, Pieter Abbeel, and Sergey Levine. A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *CoRR*, abs/1611.03852, 2016.
- [18] Ruiqi Gao, Yang Lu, Junpei Zhou, Song-Chun Zhu, and Ying Nian Wu. Learning generative convnets via multi-grid modeling and sampling. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 9155–9164, 2018.
- [19] Ruiqi Gao, Erik Nijkamp, Diederik P Kingma, Zhen Xu, Andrew M Dai, and Ying Nian Wu. Flow contrastive estimation of energy-based models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7518–7528, 2020.
- [20] Partha Ghosh, Mehdi S. M. Sajjadi, Antonio Vergari, Michael Black, and Bernhard Scholkopf. From variational to deterministic autoencoders. In *International Conference on Learning Representations*, 2020.

- [21] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010.
- [22] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014.
- [23] Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *International Conference on Learning Representations*, 2019.
- [24] Cheng-En Guo, Song-Chun Zhu, and Ying Nian Wu. Modeling visual patterns by integrating descriptive and generative methods. *International Journal of Computer Vision*, 53(1):5–29, 2003.
- [25] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, pages 297–304, 2010.
- [26] Tian Han, Yang Lu, Song-Chun Zhu, and Ying Nian Wu. Alternating back-propagation for generator network. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 1976–1984, 2017.
- [27] Tian Han, Erik Nijkamp, Xiaolin Fang, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Divergence triangle for joint training of generator model, energy-based model, and inferential model. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 8670–8679, 2019.
- [28] Tian Han, Erik Nijkamp, Xiaolin Fang, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Divergence triangle for joint training of generator model, energy-based model, and inferential model. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8670–8679, 2019.
- [29] Tian Han, Erik Nijkamp, Linqi Zhou, Bo Pang, Song-Chun Zhu, and Ying Nian Wu. Joint training of variational auto-encoder and latent energy-based model. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [30] Geoffrey E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, 2002.
- [31] Geoffrey E Hinton, Peter Dayan, Brendan J Frey, and Radford M Neal. The "wake-sleep" algorithm for unsupervised neural networks. *Science*, 268(5214):1158–1161, 1995.
- [32] Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.
- [33] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [34] Matthew D Hoffman. Learning deep latent gaussian models with markov chain monte carlo. In *International conference on machine learning*, pages 1510–1519, 2017.
- [35] Aapo Hyvärinen, Juha Karhunen, and Erkki Oja. *Independent component analysis*. John Wiley & Sons, 2004.
- [36] Long Jin, Justin Lazarow, and Zhuowen Tu. Introspective classification with convolutional nets. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 823–833, 2017.
- [37] Yoon Kim, Sam Wiseman, Andrew Miller, David Sontag, and Alexander Rush. Semi-amortized variational autoencoders. In *International Conference on Machine Learning*, pages 2678–2687, 2018.
- [38] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [39] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.

- [40] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research).
- [41] Rithesh Kumar, Anirudh Goyal, Aaron C. Courville, and Yoshua Bengio. Maximum entropy generators for energy-based models. *CoRR*, abs/1901.08508, 2019.
- [42] Justin Lazarow, Long Jin, and Zhuowen Tu. Introspective neural networks for generative modeling. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2793–2802, 2017.
- [43] Daniel D Lee and H Sebastian Seung. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*, pages 556–562, 2001.
- [44] Honglak Lee, Roger B. Grosse, Rajesh Ranganath, and Andrew Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009*, pages 609–616, 2009.
- [45] Bohan Li, Junxian He, Graham Neubig, Taylor Berg-Kirkpatrick, and Yiming Yang. A surprisingly effective fix for deep latent variable modeling of text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3603–3614, Hong Kong, China, Nov. 2019. Association for Computational Linguistics.
- [46] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [47] Yang Lu, Song-Chun Zhu, and Ying Nian Wu. Learning FRAME models using CNN filters. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pages 1902–1910, 2016.
- [48] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. Are gans created equal? a large-scale study. *arXiv preprint arXiv:1711.10337*, 2017.
- [49] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.
- [50] Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of english: The penn treebank. *Comput. Linguist.*, 19(2):313–330, June 1993.
- [51] Radford M Neal. MCMC using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, 2, 2011.
- [52] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [53] Jiquan Ngiam, Zhenghao Chen, Pang Wei Koh, and Andrew Y. Ng. Learning deep energy models. In *Proceedings of the 28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011*, pages 1105–1112, 2011.
- [54] Erik Nijkamp, Mitch Hill, Tian Han, Song-Chun Zhu, and Ying Nian Wu. On the anatomy of MCMC-based maximum likelihood learning of energy-based models. *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- [55] Erik Nijkamp, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Learning non-convergent non-persistent short-run MCMC toward energy-based model. *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, Canada*, 2019.
- [56] Erik Nijkamp, Bo Pang, Tian Han, Alex Zhou, Song-Chun Zhu, and Ying Nian Wu. Learning multi-layer latent variable model via variational optimization of short run mcmc for approximate inference. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [57] Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research*, 37(23):3311–3325, 1997.
- [58] Bo Pang, Tian Han, and Ying Nian Wu. Learning latent space energy-based prior model for molecule generation. *arXiv preprint arXiv:2010.09351*, 2020.
- [59] Bo Pang, Erik Nijkamp, Jiali Cui, Tian Han, and Ying Nian Wu. Semi-supervised learning by latent space energy-based model of symbol-vector coupling. *arXiv preprint arXiv:2010.09359*, 2020.
- [60] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.

- [61] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 1530–1538, 2015.
- [62] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pages 1278–1286, 2014.
- [63] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [64] Jason Tyler Rolfe. Discrete variational autoencoders. *arXiv preprint arXiv:1609.02200*, 2016.
- [65] Donald B. Rubin and Dorothy T. Thayer. Em algorithms for ml factor analysis. *Psychometrika*, 47(1):69–76, Mar 1982.
- [66] Ruslan Salakhutdinov and Geoffrey E. Hinton. Deep boltzmann machines. In *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics, AISTATS 2009, Clearwater Beach, Florida, USA, April 16-18, 2009*, pages 448–455, 2009.
- [67] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [68] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063, 2000.
- [69] Tijmen Tieleman. Training restricted boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th international conference on Machine learning*, pages 1064–1071, 2008.
- [70] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders. *arXiv preprint arXiv:1711.01558*, 2017.
- [71] Jakub Tomczak and Max Welling. Vae with a vampprior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223, 2018.
- [72] Jakub M. Tomczak and Max Welling. VAE with a vampprior. In Amos J. Storkey and Fernando Pérez-Cruz, editors, *International Conference on Artificial Intelligence and Statistics, AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*, volume 84 of *Proceedings of Machine Learning Research*, pages 1214–1223. PMLR, 2018.
- [73] Zhuowen Tu. Learning generative models via discriminative approaches. In *2007 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2007), 18-23 June 2007, Minneapolis, Minnesota, USA, 2007*.
- [74] Ryan D. Turner, Jane Hung, Eric Frank, Yunus Saatchi, and Jason Yosinski. Metropolis-hastings generative adversarial networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 6345–6353. PMLR, 2019.
- [75] Arash Vahdat, Evgeny Andriyash, and William Macready. Dvae#: Discrete variational autoencoders with relaxed boltzmann priors. In *Advances in Neural Information Processing Systems*, pages 1864–1874, 2018.
- [76] Arash Vahdat, William Macready, Zhengbing Bian, Amir Khoshaman, and Evgeny Andriyash. Dvae++: Discrete variational autoencoders with overlapping transformations. In *International Conference on Machine Learning*, pages 5035–5044, 2018.
- [77] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, pages 6306–6315, 2017.
- [78] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [79] Ying Nian Wu, Ruiqi Gao, Tian Han, and Song-Chun Zhu. A tale of three probabilistic families: Discriminative, descriptive, and generative models. *Quarterly of Applied Mathematics*, 77(2):423–465, 2019.
- [80] Jianwen Xie, Yang Lu, Ruiqi Gao, and Ying Nian Wu. Cooperative learning of energy-based model and latent variable model via MCMC teaching. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 4292–4301, 2018.

- [81] Jianwen Xie, Yang Lu, Song-Chun Zhu, and Ying Nian Wu. A theory of generative convnet. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 2635–2644, 2016.
- [82] Jianwen Xie, Zilong Zheng, Ruiqi Gao, Wenguan Wang, Song-Chun Zhu, and Ying Nian Wu. Learning descriptor networks for 3D shape synthesis and analysis. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 8629–8638, 2018.
- [83] Jianwen Xie, Song-Chun Zhu, and Ying Nian Wu. Synthesizing dynamic patterns by spatial-temporal generative convnet. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1061–1069, 2017.
- [84] Zichao Yang, Zhiting Hu, Ruslan Salakhutdinov, and Taylor Berg-Kirkpatrick. Improved variational autoencoders for text modeling using dilated convolutions. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 3881–3890, 2017.
- [85] Houssam Zenati, Chuan Sheng Foo, Bruno Lecouat, Gaurav Manek, and Vijay Ramaseshan Chandrasekhar. Efficient gan-based anomaly detection. *arXiv preprint arXiv:1802.06222*, 2018.
- [86] Junbo Zhao, Yoon Kim, Kelly Zhang, Alexander Rush, and Yann LeCun. Adversarially regularized autoencoders. In *International Conference on Machine Learning*, pages 5902–5911, 2018.
- [87] Song Chun Zhu. Statistical modeling and conceptualization of visual patterns. *IEEE Trans. Pattern Anal. Mach. Intell.*, 25(6):691–712, 2003.
- [88] Song Chun Zhu and David Mumford. Grade: Gibbs reaction and diffusion equations. In *Computer Vision, 1998. Sixth International Conference on*, pages 847–854, 1998.
- [89] Song Chun Zhu, Ying Nian Wu, and David Mumford. Filters, random fields and maximum entropy (FRAME): towards a unified theory for texture modeling. *International Journal of Computer Vision*, 27(2):107–126, 1998.

A Theoretical derivations

In this section, we shall derive most of the equations in the main text. We take a step by step approach, starting from simple identities or results, and gradually reaching the main results. Our derivations are unconventional, but they pertain more to our model and learning method.

A.1 A simple identity

Let $x \sim p_\theta(x)$. A useful identity is

$$\mathbb{E}_\theta[\nabla_\theta \log p_\theta(x)] = 0, \quad (18)$$

where \mathbb{E}_θ (or \mathbb{E}_{p_θ}) is the expectation with respect to p_θ .

The proof is one liner:

$$\mathbb{E}_\theta[\nabla_\theta \log p_\theta(x)] = \int [\nabla_\theta \log p_\theta(x)] p_\theta(x) dx = \int \nabla_\theta p_\theta(x) dx = \nabla_\theta \int p_\theta(x) dx = \nabla_\theta 1 = 0. \quad (19)$$

The above identity has generalized versions, such as the one underlying the policy gradient [78, 68], $\nabla_\theta \mathbb{E}_\theta[R(x)] = \mathbb{E}_\theta[R(x) \nabla_\theta \log p_\theta(x)]$. By letting $R(x) = 1$, we get (18).

A.2 Maximum likelihood estimating equation

The simple identity (18) also underlies the consistency of MLE. Suppose we observe $(x_i, i = 1, \dots, n) \sim p_{\theta_{\text{true}}}(x)$ independently, where θ_{true} is the true value of θ . The log-likelihood is

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_\theta(x_i). \quad (20)$$

The maximum likelihood estimating equation is

$$L'(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_\theta \log p_\theta(x_i) = 0. \quad (21)$$

According to the law of large number, as $n \rightarrow \infty$, the above estimating equation converges to

$$\mathbb{E}_{\theta_{\text{true}}}[\nabla_\theta \log p_\theta(x)] = 0, \quad (22)$$

where θ is the unknown value to be solved, while θ_{true} is fixed. According to the simple identity (18), $\theta = \theta_{\text{true}}$ is the solution to the above estimating equation (22), no matter what θ_{true} is. Thus with regularity conditions, such as identifiability of the model, the MLE converges to θ_{true} in probability.

The optimality of the maximum likelihood estimating equation among all the asymptotically unbiased estimating equations can be established based on a further generalization of the simple identity (18).

We shall justify our learning method with short-run MCMC in terms of an estimating equation, which is a perturbation of the maximum likelihood estimating equation (21).

A.3 MLE learning gradient for θ

Recall that $p_\theta(x, z) = p_\alpha(z) p_\beta(x|z)$, where $\theta = \{\alpha, \beta\}$. The learning gradient for an observation x is as follows:

$$\nabla_\theta \log p_\theta(x) = \mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(x, z)] = \mathbb{E}_{p_\theta(z|x)} [\nabla_\theta (\log p_\alpha(z) + \log p_\beta(x|z))]. \quad (23)$$

The above identity is a simple consequence of the simple identity (18).

$$\mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(x, z)] = \mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(z|x) + \nabla_\theta \log p_\theta(x)] \quad (24)$$

$$= \mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(z|x)] + \mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(x)] \quad (25)$$

$$= 0 + \nabla_\theta \log p_\theta(x), \quad (26)$$

because of the fact that $\mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(z|x)] = 0$ according to the simple identity (18), while $\mathbb{E}_{p_\theta(z|x)} [\nabla_\theta \log p_\theta(x)] = \nabla_\theta \log p_\theta(x)$ because what is inside the expectation only depends on x , but does not depend on z .

The above identity (23) is related to the EM algorithm [14], where x is the observed data, z is the missing data, and $\log p_\theta(x, z)$ is the complete-data log-likelihood.

A.4 MLE learning gradient for α

For the prior model $p_\alpha(z) = \frac{1}{Z(\alpha)} \exp(f_\alpha(z))p_0(z)$, we have $\log p_\alpha(z) = f_\alpha(z) - \log Z(\alpha) + \log p_0(z)$. Applying the simple identity (18), we have

$$\mathbb{E}_\alpha[\nabla_\alpha \log p_\alpha(z)] = \mathbb{E}_\alpha[\nabla_\alpha f_\alpha(z) - \nabla_\alpha \log Z(\alpha)] = \mathbb{E}_\alpha[\nabla_\alpha f_\alpha(z)] - \nabla_\alpha \log Z(\alpha) = 0. \quad (27)$$

Thus

$$\nabla_\alpha \log Z(\alpha) = \mathbb{E}_\alpha[\nabla_\alpha f_\alpha(z)]. \quad (28)$$

Hence the derivative of the log-likelihood is

$$\nabla_\alpha \log p_\alpha(x) = \nabla_\alpha f_\alpha(z) - \nabla_\alpha \log Z(\alpha) = \nabla_\alpha f_\alpha(z) - \mathbb{E}_\alpha[\nabla_\alpha f_\alpha(z)]. \quad (29)$$

According to equation (23) in the previous subsection, the learning gradient for α is

$$\nabla_\alpha \log p_\theta(x) = \mathbb{E}_{p_\theta(z|x)} [\nabla_\alpha \log p_\alpha(z)] \quad (30)$$

$$= \mathbb{E}_{p_\theta(z|x)} [\nabla_\alpha f_\alpha(z) - \mathbb{E}_{p_\alpha(z)} [\nabla_\alpha f_\alpha(z)]] \quad (31)$$

$$= \mathbb{E}_{p_\theta(z|x)} [\nabla_\alpha f_\alpha(z)] - \mathbb{E}_{p_\alpha(z)} [\nabla_\alpha f_\alpha(z)]. \quad (32)$$

A.5 Re-deriving simple identity in terms of D_{KL}

We shall provide a theoretical understanding of the learning method with short-run MCMC in terms of Kullback-Leibler divergences. We start from some simple results.

The simple identity (18) also follows from Kullback-Leibler divergence. Consider

$$D(\theta) = D_{KL}(p_{\theta_*}(x) \| p_\theta(x)), \quad (33)$$

as a function of θ with θ_* fixed. Suppose the model p_θ is identifiable, then $D(\theta)$ achieves its minimum 0 at $\theta = \theta_*$, thus $D'(\theta_*) = 0$. Meanwhile,

$$D'(\theta) = -\mathbb{E}_{\theta_*} [\nabla_\theta \log p_\theta(x)]. \quad (34)$$

Thus

$$\mathbb{E}_{\theta_*} [\nabla_\theta \log p_{\theta_*}(x)] = 0. \quad (35)$$

Since θ_* is arbitrary in the above derivation, we can replace it by a generic θ , i.e.,

$$\mathbb{E}_\theta [\nabla_\theta \log p_\theta(x)] = 0, \quad (36)$$

which is the simple identity (18).

As a notational convention, for a function $f(\theta)$, we write $f'(\theta_*) = \nabla_\theta f(\theta_*)$, i.e., the derivative of $f(\theta)$ at θ_* .

A.6 Re-deriving MLE learning gradient in terms of perturbation by D_{KL} terms

We now re-derive MLE learning gradient in terms of perturbation of log-likelihood by Kullback-Leibler divergence terms. Then the learning method with short-run MCMC can be easily understood.

At iteration t , fixing θ_t , we want to calculate the gradient of the log-likelihood function for an observation x , $\log p_\theta(x)$, at $\theta = \theta_t$. Consider the following computationally tractable perturbation of the log-likelihood

$$l_\theta(x) = \log p_\theta(x) - D_{KL}(p_{\theta_t}(z|x) \| p_\theta(z|x)) + D_{KL}(p_{\alpha_t}(z) \| p_\alpha(z)). \quad (37)$$

In the above, as a function of θ , with θ_t fixed, $D_{KL}(p_{\theta_t}(z|x) \| p_\theta(z|x))$ is minimized at $\theta = \theta_t$, thus its derivative at θ_t is 0. As a function of α , with α_t fixed, $D_{KL}(p_{\alpha_t}(z) \| p_\alpha(z))$ is minimized at $\alpha = \alpha_t$, thus its derivative at α_t is 0. Thus

$$\nabla_\theta \log p_{\theta_t}(x) = \nabla_\theta l_{\theta_t}(x). \quad (38)$$

We now unpack $l_\theta(x)$ to see that it is computationally tractable, and we can obtain its derivative at θ_t .

$$\nabla_\theta l_\theta(x) = \log p_\theta(x) + \mathbb{E}_{p_{\theta_t}(z|x)}[\log p_\theta(z|x)] - \mathbb{E}_{p_{\alpha_t}(z)}[\log p_\alpha(z)] + c \quad (39)$$

$$= \mathbb{E}_{p_{\theta_t}(z|x)}[\log p_\theta(x, z)] - \mathbb{E}_{p_{\alpha_t}(z)}[\log p_\alpha(z)] + c \quad (40)$$

$$= \mathbb{E}_{p_{\theta_t}(z|x)}[\log p_\alpha(z) + \log p_\beta(x|z)] - \mathbb{E}_{p_{\alpha_t}(z)}[\log p_\alpha(z)] + c \quad (41)$$

$$= \mathbb{E}_{p_{\theta_t}(z|x)}[\log p_\alpha(z)] - \mathbb{E}_{p_{\alpha_t}(z)}[\log p_\alpha(z)] + \mathbb{E}_{p_{\theta_t}(z|x)}[\log p_\beta(x|z)] + c \quad (42)$$

$$= \mathbb{E}_{p_{\theta_t}(z|x)}[f_\alpha(z)] - \mathbb{E}_{p_{\alpha_t}(z)}[f_\alpha(z)] + \mathbb{E}_{p_{\theta_t}(z|x)}[\log p_\beta(x|z)] + c + c', \quad (43)$$

where $\log Z(\alpha)$ term gets canceled,

$$c = -\mathbb{E}_{p_{\theta_t}(z|x)}[\log p_{\theta_t}(z|x)] + \mathbb{E}_{p_{\alpha_t}(z)}[\log p_{\alpha_t}(z)], \quad (44)$$

$$c' = \mathbb{E}_{p_{\theta_t}(z|x)}[\log p_0(z)] - \mathbb{E}_{p_{\alpha_t}(z)}[\log p_0(z)], \quad (45)$$

do not depend on θ . c consists of two entropy terms. Now taking derivative at θ_t , we have

$$\delta_{\alpha_t}(x) = \nabla_\alpha l(\theta_t) = \mathbb{E}_{p_{\theta_t}(z|x)}[\nabla_\alpha f_{\alpha_t}(z)] - \mathbb{E}_{p_{\alpha_t}(z)}[\nabla_\alpha f_{\alpha_t}(z)], \quad (46)$$

$$\delta_{\beta_t}(x) = \nabla_\beta l(\theta_t) = \mathbb{E}_{p_{\theta_t}(z|x)}[\nabla_\beta \log p_{\beta_t}(x|z)]. \quad (47)$$

Averaging over the observed examples $\{x_i, i = 1, \dots, n\}$ leads to MLE learning gradient.

In the above, we calculate the gradient of $\log p_\theta(x)$ at θ_t . Since θ_t is arbitrary in the above derivation, if we replace θ_t by a generic θ , we get the gradient of $\log p_\theta(x)$ at a generic θ , i.e.,

$$\delta_\alpha(x) = \nabla_\alpha \log p_\theta(x) = \mathbb{E}_{p_\theta(z|x)}[\nabla_\alpha f_\alpha(z)] - \mathbb{E}_{p_{\alpha_t}(z)}[\nabla_\alpha f_\alpha(z)], \quad (48)$$

$$\delta_\beta(x) = \nabla_\beta \log p_\theta(x) = \mathbb{E}_{p_\theta(z|x)}[\nabla_\beta \log p_\beta(x|z)]. \quad (49)$$

The above calculations are related to the EM algorithm [14] and the learning of energy-based model.

In EM algorithm, the complete-data log-likelihood Q serves as a surrogate for the observed-data log-likelihood $\log p_\theta(x)$, where

$$Q(\theta|\theta_t) = \log p_\theta(x) - D_{KL}(p_{\theta_t}(z|x)||p_\theta(z|x)), \quad (50)$$

and $\theta_{t+1} = \arg \max_\theta Q(\theta|\theta_t)$, where $Q(\theta|\theta_t)$ is a lower-bound of $\log p_\theta(x)$ or minorizes the latter. $Q(\theta|\theta_t)$ and $\log p_\theta(x)$ touch each other at θ_t , and they are co-tangent at θ_t . Thus the derivative of $\log p_\theta(x)$ at θ_t is the same as the derivative of $Q(\theta|\theta_t)$ at $\theta = \theta_t$.

In EBM, $D_{KL}(p_{\alpha_t}(z)||p_\alpha(z))$ serves to cancel $\log Z(\alpha)$ term in the EBM prior, and is related to the second divergence term in contrastive divergence [30].

A.7 Maximum likelihood estimating equation for $\theta = (\alpha, \beta)$

The MLE estimating equation is

$$\frac{1}{n} \sum_{i=1}^n \nabla_\theta \log p_\theta(x_i) = 0. \quad (51)$$

Based on (48) and (49), the estimating equation is

$$\frac{1}{n} \sum_{i=1}^n \delta_\alpha(x_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{p_\theta(z_i|x_i)}[\nabla_\alpha f_\alpha(z_i)] - \mathbb{E}_{p_{\alpha_t}(z)}[\nabla_\alpha f_\alpha(z)] = 0, \quad (52)$$

$$\frac{1}{n} \sum_{i=1}^n \delta_\beta(x_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{p_\theta(z_i|x_i)}[\nabla_\beta \log p_\beta(x_i|z_i)] = 0. \quad (53)$$

A.8 Learning with short-run MCMC as perturbation of log-likelihood

Based on the above derivations, we can see that learning with short-run MCMC is also a perturbation of log-likelihood, except that we replace $p_{\theta_t}(z|x)$ by $\tilde{p}_{\theta_t}(z|x)$, and replace $p_{\alpha_t}(z)$ by $\tilde{p}_{\alpha_t}(z)$, where $\tilde{p}_{\theta_t}(z|x)$ and $\tilde{p}_{\alpha_t}(z)$ are produced by short-run MCMC.

At iteration t , fixing θ_t , the updating rule based on short-run MCMC follows the gradient of the following function, which is a perturbation of log-likelihood for the observation x ,

$$\tilde{l}_\theta(x) = \log p_\theta(x) - D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_\theta(z|x)) + D_{KL}(\tilde{p}_{\alpha_t}(z)||p_\alpha(z)). \quad (54)$$

The above is a function of θ , while θ_t is fixed.

In full parallel to the above subsection, we have

$$\tilde{l}_\theta(x) = \mathbb{E}_{\tilde{p}_{\theta_t}(z|x)}[f_\alpha(z)] - \mathbb{E}_{\tilde{p}_{\alpha_t}(z)}[f_\alpha(z)] + \mathbb{E}_{\tilde{p}_{\theta_t}(z|x)}[\log p_\beta(x|z)] + c + c', \quad (55)$$

where c and c' do not depend on θ . Thus, taking derivative of the function $\tilde{l}_\theta(x)$ at $\theta = \theta_t$, we have

$$\tilde{\delta}_{\alpha_t}(x) = \nabla_\alpha \tilde{l}(\theta_t) = \mathbb{E}_{\tilde{p}_{\theta_t}(z|x)}[\nabla_\alpha f_\alpha(z)] - \mathbb{E}_{\tilde{p}_{\alpha_t}(z)}[\nabla_\alpha f_\alpha(z)], \quad (56)$$

$$\tilde{\delta}_{\beta_t}(x) = \nabla_\beta \tilde{l}(\theta_t) = \mathbb{E}_{\tilde{p}_{\theta_t}(z|x)}[\nabla_\beta \log p_\beta(x|z)]. \quad (57)$$

Averaging over $\{x_i, i = 1, \dots, n\}$, we get the updating rule based on short-run MCMC. That is, the learning rule based on short-run MCMC follows the gradient of a perturbation of the log-likelihood function where the perturbations consists of two D_{KL} terms.

$D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_\theta(z|x))$ is related to VAE [39], where $\tilde{p}_{\theta_t}(z|x)$ serves as an inference model, except that we do not learn a separate inference network. $D_{KL}(\tilde{p}_{\alpha_t}(z)||p_\alpha(z))$ is related to contrastive divergence [30], except that $\tilde{p}_{\alpha_t}(z)$ is initialized from the Gaussian white noise $p_0(z)$, instead of the data distribution of observed examples.

$D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_\theta(z|x))$ and $D_{KL}(\tilde{p}_{\alpha_t}(z)||p_\alpha(z))$ cause the bias relative to MLE learning. MLE is impractical because we cannot do exact sampling with MCMC.

However, the bias may not be all that bad. In learning β , $D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_\theta(z|x))$ may force the model to be biased towards the approximate short-run posterior $\tilde{p}_{\theta_t}(z|x)$, so that the short-run posterior is close to the true posterior. In learning α , the update based on $\mathbb{E}_{\tilde{p}_{\theta_t}(z|x)}[\nabla_\alpha f_\alpha(z)] - \mathbb{E}_{\tilde{p}_{\alpha_t}(z)}[\nabla_\alpha f_\alpha(z)]$ may force the short-run prior $\tilde{p}_\alpha(z)$ to match the short-run posterior $\tilde{p}_\theta(z|x)$.

A.9 Perturbation of maximum likelihood estimating equation

The fixed point of the learning algorithm based on short-run MCMC is where the update is 0, i.e.,

$$\frac{1}{n} \sum_{i=1}^n \tilde{\delta}_\alpha(x_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{p}_\theta(z_i|x_i)}[\nabla_\alpha f_\alpha(z_i)] - \mathbb{E}_{\tilde{p}_\alpha(z)}[\nabla_\alpha f_\alpha(z)] = 0, \quad (58)$$

$$\frac{1}{n} \sum_{i=1}^n \tilde{\delta}_\beta(x_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{p}_\theta(z_i|x_i)}[\nabla_\beta \log p_\beta(x_i|z_i)] = 0. \quad (59)$$

This is clearly a perturbation of the MLE estimating equation in (52) and (53). The above estimating equation defines an estimator, where the learning algorithm with short-run MCMC converges.

A.10 Three D_{KL} terms

We can rewrite the objective function (54) in a more revealing form. Let $(x_i, i = 1, \dots, n) \sim p_{\text{data}}(x)$ independently, where $p_{\text{data}}(x)$ is the data distribution. At time step t , with fixed θ_t , learning based on short-run MCMC follows the gradient of

$$\frac{1}{n} \sum_{i=1}^n [\log p_\theta(x_i) - D_{KL}(\tilde{p}_{\theta_t}(z_i|x_i)||p_\theta(z_i|x_i)) + D_{KL}(\tilde{p}_{\alpha_t}(z)||p_\alpha(z))]. \quad (60)$$

Let us assume n is large enough, so that the average is practically the expectation with respect to p_{data} . Then MLE maximizes $\frac{1}{n} \sum_{i=1}^n \log p_\theta(x_i) \doteq \mathbb{E}_{p_{\text{data}}(x)}[\log p_\theta(x)]$, which is equivalent to minimizing $D_{KL}(p_{\text{data}}(x)||p_\theta(x))$. The learning with short-run MCMC follows the gradient that minimizes

$$D_{KL}(p_{\text{data}}(x)||p_\theta(x)) + D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_\theta(z|x)) - D_{KL}(\tilde{p}_{\alpha_t}(z)||p_\alpha(z)), \quad (61)$$

where, with some abuse of notation, we now define

$$D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_\theta(z|x)) = \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{\tilde{p}_{\theta_t}(z|x)} \left[\log \frac{\tilde{p}_{\theta_t}(z|x)}{p_\theta(z|x)} \right], \quad (62)$$

where we also average over $x \sim p_{\text{data}}(x)$, instead fixing x as before.

The objective (61) is clearly a perturbation of the MLE, as the MLE is based on the first D_{KL} in (61). The signs in front of the remaining two D_{KL} perturbations also become clear. The sign in front of $D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_{\theta}(z|x))$ is positive because

$$D_{KL}(p_{\text{data}}(x)||p_{\theta}(x)) + D_{KL}(\tilde{p}_{\theta_t}(z|x)||p_{\theta}(z|x)) = D_{KL}(p_{\text{data}}(x)\tilde{p}_{\theta_t}(z|x)||p_{\alpha}(x)p_{\beta}(x|z)), \quad (63)$$

where the D_{KL} on the right hand side is about the joint distributions of (x, z) , and is more tractable than the first D_{KL} on the left hand side, which is for MLE. This underlies EM and VAE. Now subtracting the third D_{KL} , we have the following special form of contrastive divergence

$$D_{KL}(p_{\text{data}}(x)\tilde{p}_{\theta_t}(z|x)||p_{\alpha}(z)p_{\beta}(x|z)) - D_{KL}(\tilde{p}_{\alpha_t}(z)||p_{\alpha}(z)), \quad (64)$$

where the negative sign in front of $D_{KL}(\tilde{p}_{\alpha_t}(z)||p_{\alpha}(z))$ is to cancel the intractable $\log Z(\alpha)$ term.

The above contrastive divergence also has an adversarial interpretation. When $p_{\alpha}(z)$ or α is updated, $p_{\alpha}(z)p_{\beta}(x|z)$ gets closer to $p_{\text{data}}(x)\tilde{p}_{\theta_t}(z|x)$, while getting away from $\tilde{p}_{\alpha_t}(z)$, i.e., p_{α} seeks to criticize the samples from $\tilde{p}_{\alpha_t}(z)$ by comparing them to the posterior samples of z inferred from the real data.

As mentioned in the main text, we can also exponentially tilt $p_0(x, z) = p_0(z)p_{\beta}(x|z)$ to $p_{\theta}(x, z) = \frac{1}{Z(\theta)} \exp(f_{\alpha}(x, z))p_0(x, z)$, or equivalently, exponentially tilt $p_0(z, \epsilon) = p_0(z)p(\epsilon)$. The above derivations can be easily adapted to such a model, which we choose not to explore due to the complexity of EBM in the data space.

A.11 Amortized inference and synthesis networks

We can jointly train two extra networks together with the original model to amortize the short-run MCMC for inference and synthesis sampling. Specifically, we use an inference network $q_{\phi}(z|x)$ to amortize the short-run MCMC that produces $\tilde{p}_{\theta}(z|x)$, and we use a synthesis network $q_{\psi}(z)$ to amortize the short-run MCMC that produces $\tilde{p}_{\alpha}(z)$.

We can then define the following objective function in parallel with the objective function (61) in the above subsection,

$$\Delta(\theta, \phi, \psi) = D_{KL}(p_{\text{data}}(x)||p_{\theta}(x)) + D_{KL}(q_{\phi}(z|x)||p_{\theta}(z|x)) - D_{KL}(q_{\psi}(z)||p_{\alpha}(z)), \quad (65)$$

and we can jointly learn θ , ϕ and ψ by

$$\min_{\theta} \min_{\phi} \max_{\psi} \Delta(\theta, \phi, \psi). \quad (66)$$

See [28, 29] for related formulations. The learning of the inference network $q_{\phi}(z|x)$ follows VAE. The learning of the synthesis network $q_{\psi}(z)$ is based on variational approximation to $p_{\alpha}(z)$. The pair $p_{\alpha}(z)$ and $q_{\psi}(z)$ play adversarial roles, where $q_{\psi}(z)$ serves as an actor and $p_{\alpha}(z)$ serves as a critic.

B Experiments

B.1 Experiment details

Data. Image datasets include SVHN [52] ($32 \times 32 \times 3$), CIFAR-10 [40] ($32 \times 32 \times 3$), and CelebA [46] ($64 \times 64 \times 3$). We use the full training split of SVHN (73, 257) and CIFAR-10 (50, 000) and take 40, 000 examples of CelebA as training data following [55]. The training images are resized and scaled to $[-1, 1]$. Text datasets include PTB [50], Yahoo [84], and SNLI [5], following recent work on text generative modeling with latent variables [37, 86, 45].

Model architectures. The architecture of the EBM, $f_{\alpha}(z)$, is displayed in Table 6. For text data, the dimensionality of z is set to 32. The generator architectures for the image data are also shown in Table 6. The generators for the text data are implemented with a one-layer unidirectional LSTM [33] and Table 7 lists the number of word embeddings and hidden units of the generators for each dataset.

Short run dynamics. The hyperparameters for the short run dynamics are depicted in Table 5 where K_0 and K_1 denote the number of prior and posterior sampling steps with step sizes s_0 and

s_1 , respectively. These are identical across models and data modalities, except for the model for CIFAR-10 which is using $K_1 = 40$ steps.

Short Run Dynamics Hyperparameters	
Hyperparameter	Value
K_0	60
s_0	0.4
K_1	20
s_1	0.1

Table 5: Hyperparameters for short run dynamics.

Optimization. The parameters for the EBM and image generators are initialized with Xavier normal [21] and those for the text generators are initialized from a uniform distribution, $\text{Unif}(-0.1, 0.1)$, following [37, 45]. Adam [38] is adopted for all model optimization. The models are trained until convergence (taking approximately 70,000 and 40,000 parameter updates for image and text models, respectively).

EBM Model		
Layers	In-Out Size	Stride
Input: z	100	
Linear, LReLU	200	-
Linear, LReLU	200	-
Linear	1	-
Generator Model for SVHN, ngf = 64		
Input: x	1x1x100	
4x4 convT(ngf x 8), LReLU	4x4x(ngf x 8)	1
4x4 convT(ngf x 4), LReLU	8x8x(ngf x 4)	2
4x4 convT(ngf x 2), LReLU	16x16x(ngf x 2)	2
4x4 convT(3), Tanh	32x32x3	2
Generator Model for CIFAR-10, ngf = 128		
Input: x	1x1x128	
8x8 convT(ngf x 8), LReLU	8x8x(ngf x 8)	1
4x4 convT(ngf x 4), LReLU	16x16x(ngf x 4)	2
4x4 convT(ngf x 2), LReLU	32x32x(ngf x 2)	2
3x3 convT(3), Tanh	32x32x3	1
Generator Model for CelebA, ngf = 128		
Input: x	1x1x100	
4x4 convT(ngf x 8), LReLU	4x4x(ngf x 8)	1
4x4 convT(ngf x 4), LReLU	8x8x(ngf x 4)	2
4x4 convT(ngf x 2), LReLU	16x16x(ngf x 2)	2
4x4 convT(ngf x 1), LReLU	32x32x(ngf x 1)	2
4x4 convT(3), Tanh	64x64x3	2

Table 6: EBM model architectures for all image and text datasets and generator model architectures for SVHN ($32 \times 32 \times 3$), CIFAR-10 ($32 \times 32 \times 3$), and CelebA ($64 \times 64 \times 3$). convT(n) indicates a transposed convolutional operation with n output feature maps. LReLU indicates the Leaky-ReLU activation function. The leak factor for LReLU is 0.2 in EBM and 0.1 in Generator.

	SNLI	PTB	Yahoo
Word Embedding Size	256	128	512
Hidden Size of Generator	256	512	1024

Table 7: The sizes of word embeddings and hidden units of the generators for SNLI, PTB, and Yahoo.

C Ablation study

We investigate a range of factors that are potentially affecting the model performance with SVHN as an example. The highlighted number in Tables 8, 9, and 10 is the FID score reported in the main text and compared to other baseline models. It is obtained from the model with the architecture and hyperparameters specified in Table 5 and Table 6 which serve as the reference configuration for the ablation study.

Fixed prior. We examine the expressivity endowed with the EBM prior by comparing it to models with a fixed isotropic Gaussian prior. The results are displayed in Table 8. The model with an EBM prior clearly outperforms the model with a fixed Gaussian prior and the same generator as the reference model. The fixed Gaussian models exhibit an enhancement in performance as the generator complexity increases. They however still have an inferior performance compared to the model with an EBM prior even when the fixed Gaussian prior model has a generator with four times more parameters than that of the reference model.

Model	FID
Latent EBM Prior	29.44
Fixed Gaussian	
same generator	43.39
generator with 2 times as many parameters	41.10
generator with 4 times as many parameters	39.50

Table 8: Comparison of the models with a latent EBM prior versus a fixed Gaussian prior. The highlighted number is the reported FID for SVHN and compared to other baseline models in the main text.

MCMC steps. We also study how the number of short run MCMC steps for prior inference (K_0) and posterior inference (K_1). The left panel of Table 9 shows the results for K_0 and the right panel for K_1 . As the number of MCMC steps increases, we observe improved quality of synthesis in terms of FID.

Steps	FID	Steps	FID
$K_0 = 40$	31.49	$K_1 = 20$	29.44
$K_0 = 60$	29.44	$K_1 = 40$	27.26
$K_0 = 80$	28.32	$K_1 = 60$	26.13

Table 9: Influence of the number of prior and posterior short run steps K_0 (left) and K_1 (right). The highlighted number is the reported FID for SVHN and compared to other baseline models in the main text.

Prior EBM and generator complexity. Table 10 displays the FID scores as a function of the number of hidden features of the prior EBM (**nef**) and the factor of the number of channels of the generator (**ngf**, also see Table 6). In general, enhanced model complexity leads to improved generation.

nef	50	100	200
32	32.25	31.98	30.78
ngf 64	30.91	30.56	29.44
128	29.12	27.24	26.95

Table 10: Influence of prior and generator complexity. The highlighted number is the reported FID for SVHN and compared to other baseline models in the main text. **nef** indicates the number of hidden features of the prior EBM and **ngf** denotes the factor of the number of channels of the generator (also see Table 6).

D PyTorch code

```
import torch as t, torch.nn as nn
import torchvision as tv, torchvision.transforms as tfm

img_size, batch_size = 32, 100
nz, nc, ndf, ngf = 100, 3, 200, 64
K_0, a_0, K_1, a_1 = 60, 0.4, 40, 0.1
llhd_sigma = 0.3
n_iter = 70000
device = t.device('cuda' if t.cuda.is_available() else 'cpu')

class _G(nn.Module):
    def __init__(self):
        super().__init__()
        self.gen = nn.Sequential(nn.ConvTranspose2d(nz, ngf*8, 4, 1, 0), nn.LeakyReLU(),
                                   nn.ConvTranspose2d(ngf*8, ngf*4, 4, 2, 1), nn.LeakyReLU(),
                                   nn.ConvTranspose2d(ngf*4, ngf*2, 4, 2, 1), nn.LeakyReLU(),
                                   nn.ConvTranspose2d(ngf*2, nc, 4, 2, 1), nn.Tanh())
    def forward(self, z):
        return self.gen(z)

class _E(nn.Module):
    def __init__(self):
        super().__init__()
        self.ebm = nn.Sequential(nn.Linear(nz, ndf), nn.LeakyReLU(0.2),
                                   nn.Linear(ndf, ndf), nn.LeakyReLU(0.2),
                                   nn.Linear(ndf, 1))
    def forward(self, z):
        return self.ebm(z.squeeze()).view(-1, 1, 1, 1)

transform = tfm.Compose([tfm.Resize(img_size), tfm.ToTensor(), tfm.Normalize([(0.5]*3), ([0.5]*3)],)]
data = t.stack([x[0] for x in tv.datasets.SVHN(root='data/svhn', transform=transform)]).to(device)

G, E = _G().to(device), _E().to(device)
mse = nn.MSELoss(reduction='sum')
optE = t.optim.Adam(E.parameters(), lr=0.00002, betas=(0.5, 0.999))
optG = t.optim.Adam(G.parameters(), lr=0.0001, betas=(0.5, 0.999))

def sample_p_data():
    return data[t.LongTensor(batch_size).random_(0, data.size(0))].detach()

def sample_p_0(n=batch_size):
    return t.randn(*[n, nz, 1, 1]).to(device)

def sample_langevin_prior(z, E):
    z = z.clone().detach().requires_grad_(True)
    for i in range(K_0):
        en = E(z)
        z_grad = t.autograd.grad(en.sum(), z)[0]
        z.data = z.data - 0.5 * a_0 * a_0 * (z_grad + 1.0 / z.data) + a_0 * t.randn_like(z).data
    return z.detach()

def sample_langevin_posterior(z, x, G, E):
    z = z.clone().detach().requires_grad_(True)
    for i in range(K_1):
        x_hat = G(z)
        g_log_lkhd = 1.0 / (2.0 * llhd_sigma * llhd_sigma) * mse(x_hat, x)
        grad_g = t.autograd.grad(g_log_lkhd, z)[0]
        en = E(z)
        grad_e = t.autograd.grad(en.sum(), z)[0]
        z.data = z.data - 0.5 * a_1 * a_1 * (grad_g + grad_e + 1.0 / z.data) + a_1 * t.randn_like(z).data
    return z.detach()

for i in range(n_iter):
    x = sample_p_data()
    z_e_0, z_g_0 = sample_p_0(), sample_p_0()
    z_e_k, z_g_k = sample_langevin_prior(z_e_0, E), sample_langevin_posterior(z_g_0, x, G, E)

    optG.zero_grad()
    x_hat = G(z_g_k.detach())
    loss_g = mse(x_hat, x) / batch_size
    loss_g.backward()
    optG.step()

    optE.zero_grad()
    en_pos, en_neg = E(z_g_k.detach()).mean(), E(z_e_k.detach()).mean()
    loss_e = en_pos - en_neg
    loss_e.backward()
    optE.step()
```