

# LATENT-INFORMED ENERGY-BASED MODELS WITH COLLABORATIVE GENERATOR TRAINING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Energy-based models (EBMs) have established a distinct niche in generative modeling through their architectural flexibility and expressive density estimation. However, they have yet to achieve mainstream adoption due to their training challenges. In this paper, we propose training latent-informed EBMs that leverage self-supervised representation learning to derive informative target latent variables. This joint space optimization enables the energy function to capture both data distribution and semantic manifold geometry. To avoid long-run MCMC sampling, we introduce an auxiliary generator with specialized training designs for effective energy-generator collaboration. Our training paradigm only requires MCMC sampling in the data space, and the joint energy function learns semantic data-latent relationships directly from real data. Experiments show our approach significantly boosts the generation performance compared to current EBMs with fewer MCMC steps and smaller networks. We also demonstrate the capabilities of our model across multiple tasks, including out-of-distribution detection, conditional sampling, and zero-shot image restoration.

## 1 INTRODUCTION

Generative models have achieved unprecedented rapid development in recent years. Energy-based models (EBMs) (LeCun et al., 2006; Salakhutdinov et al., 2007; Du & Mordatch, 2019), as a class of generative models, occupy a unique position among various generative frameworks due to their huge potential in modeling complex data distributions. With a flexible energy function to directly characterize the underlying probability distribution, EBM can be useful in various tasks such as image and video synthesis (Xie et al., 2019; Zhao et al., 2020), image restoration (Xie et al., 2021a; Gao et al., 2021), compositional generation (Du et al., 2020; 2023), and out-of-distribution (OOD) detection (Yoon et al., 2023; 2021). However, it is notorious for hard training and long-run MCMC sampling (Grathwohl et al., 2021; Nijkamp et al., 2020), leaving a noticeable gap with dominant generative models.

Adversarial EBMs (Geng et al., 2021; 2024) and cooperative learning (Xie et al., 2020) incorporate a generator to speed up sampling and improve generation quality. However, adversarial EBMs are prone to suffering from mode collapse because of their minimax training strategy. Cooperative learning leads to biased generator learning, thereby limiting the potential for learning a robust EBM. Divergence Triangle methods (Han et al., 2019; 2020) extend this co-training scheme to latent-variable models. However, by enforcing exact alignment between the latent representation and the generator’s prior distribution, they restrict both generation quality and latent space flexibility, ultimately weakening the energy function. CLEL (Lee et al., 2023) designs a new class of latent-variable EBMs that model the joint distribution using a contrastive latent encoder. This architecture enables the energy function to benefit from the semantically informative latent representation, moving beyond the conventional Gaussian posterior. But its training paradigm is inelegant, and sampling remains slow.

We propose a collaborative training framework that combines latent-informed EBMs (LIEBMs) with auxiliary generator initialization. For each training step, the energy function and generator are updated alternatively. When training LIEBM, the defined energy distribution is optimized in a joint space, where the target latent variables are derived through a pretrained self-supervised latent encoder. This design helps energy function understand the semantic geometry of the data manifold, as shown in Fig.1. Samples from the energy distribution are required for training as negative

054 samples. We obtain these via generator-predicted initial samples, followed by brief MCMC sam-  
 055 pling. An augmentation technique is applied to negative samples to improve the energy function’s  
 056 discrimination of regions that are far away from the data distribution. Our generator learns to ap-  
 057 proximate the long-run MCMC dynamics through a single-step transformation, thereby enabling  
 058 efficient short-run sampling that avoids the slow convergence of traditional MCMC approaches. We  
 059 investigate several designs of generator learning and conduct a thorough comparison among them.  
 060 Our method improves EBM performance through dedicated collaborative training framework and  
 061 isolates the generative prior from semantic latent representation, thus avoiding the potential pitfalls  
 062 of “posterior collapse” (Geng et al., 2023).

063 Our main contributions are summarized as follows:

- 064 1. We introduce a unified and efficient latent-informed EBM, demonstrating the necessity of  
 065 joint space optimization. Our training paradigm requires MCMC sampling only in the data  
 066 space, allowing the joint energy function to learn semantic data–latent relationships directly  
 067 from real samples, rather than between real and synthetic spaces.
- 068 2. Our approach utilizes pretrained self-supervised representations as latent variables, main-  
 069 taining their independence from the generator’s prior. This architecture provides semantic  
 070 guidance to the energy function to tap into full generative potential.
- 071 3. We develop a collaborative training framework between the EBM and an auxiliary gener-  
 072 ator, incorporating key design choices that enable effective collaboration, including a  
 073 negative-sample augmentation strategy and adaptive generator training schemes.
- 074 4. Our method achieves superior sample quality with lightweight architectures, while exhibit-  
 075 ing versatile applicability across multiple downstream tasks, including OOD detection,  
 076 conditional sampling, and zero-shot image restoration.

## 079 2 PRELIMINARY

080 Latent-variable EBMs generalize standard EBMs by incorporating a latent variable to model a joint  
 081 distribution  $p_\theta(x, z)$ :

$$082 \quad p_\theta(x, z) = \frac{\exp(E_\theta(x, z))}{Z_\theta}, \quad Z_\theta = \int \exp(E_\theta(x, z)) dx dz, \quad (1)$$

083 where  $Z_\theta$  is the intractable normalizing constant called the partition function. Training latent-  
 084 variable EBMs primarily relies on maximizing the log-likelihood such that:

$$085 \quad L(\theta) := \mathbb{E}_{(x,z) \sim p_{\text{data}}(x,z)} [\log p_\theta(x, z)] = \mathbb{E}_{(x,z) \sim p_{\text{data}}(x,z)} [E_\theta(x, z)] - \log Z_\theta. \quad (2)$$

086 Similar to standard EBMs, the gradient of the training objective can be written as:

$$087 \quad \frac{\partial L}{\partial \theta} = \mathbb{E}_{(x,z) \sim p_{\text{data}}(x,z)} \left[ \frac{\partial}{\partial \theta} E_\theta(x, z) \right] - \mathbb{E}_{(x,z) \sim p_\theta(x,z)} \left[ \frac{\partial}{\partial \theta} E_\theta(x, z) \right]. \quad (3)$$

088 It requires MCMC sampling from a joint distribution  $p_\theta(x, z)$ , which can be challenging in complex  
 089 high-dimensional space (Xu et al., 2018). Alternatively, Eq.3 can be reformulated to require only  
 090 sampling from the marginal distribution  $p_\theta(x)$ :

$$091 \quad \frac{\partial L}{\partial \theta} = \mathbb{E}_{(x,z) \sim p_{\text{data}}(x,z)} \left[ \frac{\partial}{\partial \theta} E_\theta(x, z) \right] - \mathbb{E}_{x \sim p_\theta(x)} \left[ \frac{\partial}{\partial \theta} E_\theta(x) \right], \quad (4)$$

092 where  $E_\theta(x) = \log \int \exp(E_\theta(x, z)) dz$  is an available energy function of marginal  $p_\theta(x)$ . See  
 093 Appendix A.3 and A.4 for derivation and additional details about EBMs.

## 094 3 METHOD

095 Conventional EBMs train the energy function solely in the data space, which poses challenges in  
 096 high-dimensional settings due to data sparsity and limited distributional information. To address  
 097 this, we propose a latent-variable EBM with structured latent constraints and generator-assisted

MCMC initialization. The energy function and generator are trained alternatively within each training step. Our formulation needs to solve three fundamental problems: defining a target joint distribution  $p_{\text{data}}(x, z)$  given only observed samples  $x$ , constructing a joint energy distribution  $p_{\theta}(x, z)$  that captures data-latent coupling, and balancing collaborative training between the energy function and generator.

### 3.1 LATENT-INFORMED EBM TRAINING

We define a conditional latent distribution by mapping data samples to latent variables through a latent encoder, i.e., sampling  $p_{\text{data}}(z|x)$  through  $h(v(x))/\|h(v(x))\|_2, v \sim \mathcal{V}$ , where  $\mathcal{V}$  is a random operation distribution. Our latent encoder  $h$  is pretrained using self-supervised representation learning as a separate stage before EBM training. We observe in Fig.7 that a CLEL-style collaborative training approach leads to measurable degradation in the encoder’s classification accuracy, which subsequently impairs EBM training. This phenomenon may stem from our generator-initialized EBM samples inadequately covering the true data manifold in the early stage of training, making their latent variables ineffective as negative representations for diversity.

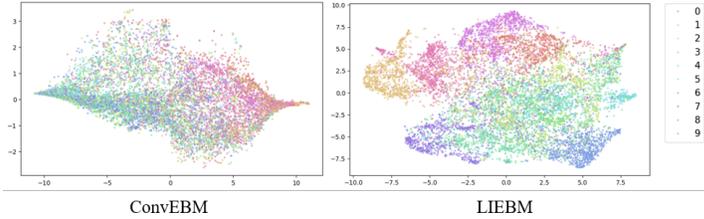


Figure 1: t-SNE visualization of  $f_{\phi}(x)$  trained on CIFAR-10: conventional EBM with  $F(f_{\phi}(x))$  as energy function vs. our LIEBM.

Considering modeling a joint energy distribution, we define our energy function by decomposing the joint density into an implicit data distribution and an explicit latent posterior:

$$E_{\theta}(x, z) = F(f_{\phi}(x)) + \log p_{\phi, \psi}(z|x), \quad (5)$$

$$E_{\theta}(x) = \log \int \exp(E_{\theta}(x, z)) dz = F(f_{\phi}(x)), \quad (6)$$

where  $f_{\phi}(x)$  is a neural network parameterized by  $\phi$ ,  $F$  maps  $f_{\phi}(x)$  to a scalar value, which can be a non-parametric function or a neural network.  $p_{\phi, \psi}(z|x)$  is a probability density parameterized by  $(\phi, \psi)$ , and  $\theta = (\phi, \psi)$ . This formulation permits EBM training via Eq.4, requiring only that  $p_{\phi, \psi}(z|x)$  be an explicit density function. Employing Eq.4 as training objective only requires MCMC sampling in data space, avoiding expensive cost in augmented  $(x, z)$  space. Moreover, via Eq.4, our joint energy function learns semantic data-latent relationships directly from real data, which is more reasonable than previous contrastive learning between real and fake spaces in EBMs and GANs.

Our joint energy definition in Eq.5 is a general formulation that theoretically encompasses most existing latent-variable EBM variants. When  $p_{\phi, \psi}(z|x)$  is defined as Gaussian, our  $E_{\theta}(x, z)$  reduces to conventional formulations (Cui & Han, 2023; Han et al., 2020; Kan et al., 2022); With cosine-similarity form, it’s similar to CLEL; Furthermore, with an exponential family form of  $E_{\theta}(x, z) = \langle \lambda - f_{\phi}(x), \eta(z) \rangle + B(\lambda)$ , our definition reduces to CEEM (Wu et al., 2021). In this case,  $E_{\theta}(x) = B(\lambda + f_{\phi}(x)) - B(\lambda)$  and  $p_{\phi, \psi}(z|x) = \langle \eta(z), \lambda + f_{\phi}(x) \rangle - B(\lambda + f_{\phi}(x))$ . With  $\lambda = (\lambda_1, \lambda_2) = (0, -\frac{1}{2})$ ,  $B(\lambda) = -\frac{\lambda_1^2}{4\lambda_2} - \frac{1}{2} \log(-2\lambda_2)$ , and  $\eta(z_k) = (z_k, z_k^2)$  for each dimension of  $z$ ,  $p_{\phi, \psi}(z|x)$  becomes a Gaussian distribution. We empirically compare these choices, and the cosine-similarity is stable and performs significantly better than the others (See Sec.4.5). Therefore, we adopt cosine-similarity form to define posterior  $p_{\phi, \psi}(z|x)$  on a unit sphere:

$$p_{\phi, \psi}(z|x) = \frac{\exp(\gamma \text{sim}(g_{\psi}(f_{\phi}(x)), z))}{Z_{\gamma}}, \quad z \sim \mathbb{S}^{d_z-1} \quad (7)$$

where  $\text{sim}(u, v) = u^{\top} v / (\|u\|_2 \|v\|_2)$  is the cosine similarity. This definition revisits the conventional use of Gaussian latents with benefits from two critical properties: (1) the normalizing constant  $Z_{\gamma}$  is independent of  $\theta$ , which can be omitted during training; (2) the scale hyperparameter  $\gamma$  controls density magnitudes for training stability.

We optimize our energy function using Eq.4, which requires sampling negative samples from the marginal  $p_{\theta}(x)$ . To avoid long MCMC chains, we consider first generating initial samples through a

generator, i.e.,  $x^0 = G(m)$ ,  $m \sim \mathcal{N}(0, I)$ , then refining them with a few MCMC steps from  $E_\theta(x^t)$ . However, this strategy exhibits a practical limitation: the initial sample distribution progressively becomes closer to the data distribution during training, resulting in the energy function’s catastrophic forgetting of low-density regions and earlier discovered modes. To mitigate this problem, we implement a **stochastic augmentation strategy** for negative samples before MCMC sampling. Each negative sample undergoes augmentation with Bernoulli probability  $p$ , where the augmented transformation  $v \sim \mathcal{V}$  follows the same protocol as used for sampling from  $p_{\text{data}}(z|x)$ . This augmentation technique enables broader exploration of the energy landscape during training, facilitating diversity of MCMC chains. Empirically, this augmentation technique enhances OOD detection for distant outliers with minimal impact on generation quality.

### 3.2 GENERATOR TRAINING

We introduce a generator to initialize MCMC chains via single-step forward propagation. This generator is typically optimized through adversarial training or cooperative learning. We build on cooperative learning as adversarial training would necessitate computationally challenging entropy maximization of the generated distribution  $p_g(x)$ . Beyond cooperative learning, our framework features a joint energy function and a semantic-aware latent encoder. These architectural advantages allow us to investigate distinct generator training schemes through extensive empirical analysis.

#### 3.2.1 ENERGY DISTRIBUTION MATCHING (EM)

Following cooperative learning, the generator can be optimized by minimizing the KL divergence between two joint distributions,  $\min \text{KL}(p_\theta(x, m) \| p_g(x, m))$ , both distributions built from a Gaussian prior  $p(m)$  and conditional  $p(x|m)$ . Under the assumption that  $p_g(x|m)$  follows a Gaussian distribution, this objective simplifies to an MSE loss:

$$L_G = \frac{1}{n} \sum_{i=1}^n \frac{1}{2\tau^2} \|G(m_i) - x_i^T\|_2^2, \quad (8)$$

where  $i$  denotes  $i^{\text{th}}$  number of a batch with size  $n$ .  $\tau^2$  is the fixed variance of  $p_g(x|m)$ .  $x_i^T$  is the refined samples by running  $T$  steps of MCMC from initial point  $x_i^0 = G(m_i)$ :

$$x_i^{t+1} = x_i^t + \frac{\delta^2}{2} \nabla_x E_\theta(x_i^t) + \delta \epsilon^t, \quad \epsilon^t \sim \mathcal{N}(0, I) \quad (9)$$

$E_\theta$  is the marginal energy defined in Eq.6. We empirically observe that this marginal energy MCMC performs well across all evaluated datasets.

We also investigate MCMC refinement from the perspective of the joint energy function. We take the basic idea of auxiliary variable MCMC (Brooks et al., 2011; Song & Ou, 2018) to sample in the augmented space  $(x, z)$ . To circumvent the computational burden of two Markov chains in both data and latent spaces, we employ our latent encoder to perform a single MCMC procedure. Specifically, we first sample initial  $x_i^0 = G(m)$ , followed by executing MCMC as described below:

$$x_i^{t+1} = x_i^t + \frac{\delta^2}{2} \nabla_x E_\theta(x_i^t, h(x_i^t)) + \delta \epsilon^t, \quad \epsilon^t \sim \mathcal{N}(0, I) \quad (10)$$

This approach is justified because  $p_\theta(z|x)$  is learned to match  $p_{\text{data}}(z|x)$  during EBM training, while constraining  $x$  within the latent space reduces the search space and improves efficiency. We observe that this joint energy refinement accelerates training in the early stage, but ultimately underperforms marginal MCMC when dealing with multimodal distributions.

#### 3.2.2 ENERGY AND REAL DISTRIBUTION MATCHING (ERM)

Dual-MCMC (Cui & Han, 2023) highlights that energy distribution matching may induce biased generator learning because it solely aligns with the energy distribution without direct access to training data. Inspired by Dual-MCMC, also leveraging our latent encoder, we optimize the generator to match both the energy and real data distribution, yielding a more informative initialization.

$$L_G = \omega_1 \text{KL}(p_\theta(x, m) \| p_g(x, m)) + \omega_2 \text{KL}(p_{\text{data}}(x, z) \| p_g(x, z)), \quad (11)$$

where  $\omega_1$  and  $\omega_2$  denote the importance weighting between two divergence components. The first term is equal to Eq.8. For the second term, we define  $p_g(x, z)$  as:

$$p_g(x, z) = \int p(m)p_g(x, z|m)dm, \quad (12)$$

$$\log p_g(x, z|m) = \log p_g(x|m) + \rho \text{sim}(z, h(G(m))) - \log Z_\rho. \quad (13)$$

We define  $p_g(x, z)$  in this way to facilitate both latent alignment and pixel-level fidelity. The second term in Eq.11 can be optimized by the classic evidence lower bound (ELBO):

$$-\mathbb{E}_{p_{\text{data}}(x,z)}\mathbb{E}_{q_\alpha(m|x)} \left[ \log p_g(x, z|m) - \log \frac{q_\alpha(m|x)}{p(m)} \right] \quad (14)$$

where  $q_\alpha$  denotes an inference model parameterized by  $\alpha$  and jointly trained with the generator. This method introduces an extra network, while increasing training complexity, this autoencoder-based architecture would be necessary for applications such as image restoration.

Table 1: Generative performance on CIFAR-10. “w/o MCMC” denotes direct sampling from the generator without energy-based refinement via MCMC sampling.

Model	NFE↓	FID↓	IS↑	Model	NFE↓	FID↓	IS↑
<b>Likelihood-based</b>				<b>EBM-based</b>			
PixelCNN (Oord et al., 2016)	1024	65.9	4.60	IGEBM (Du & Mordatch, 2019)	60	38.2	6.78
Glow (Kingma & Dhariwal, 2018)	1	48.9	3.92	joint Triangle (Han et al., 2020)	1	30.10	7.17
VAE (Kingma & Welling, 2014)	1	115.8	3.8	CoopNets (Xie et al., 2020)	51	33.61	6.55
NVAE (Vahdat & Kautz, 2020)	1	51.67	5.51	EBMBB (Geng et al., 2021)	1	28.63	7.45
<b>GAN-based</b>				VAEBM (Xiao et al., 2021)	16	12.19	8.43
SN-GAN (Miyato et al., 2018)	1	21.7	8.22	DRL (Gao et al., 2021)	180	9.58	8.30
BigGAN (Brock et al., 2019)	1	14.73	9.22	CoopFlow (Xie et al., 2022)	31	15.80	–
StyleGAN2 w/ ADA(Karras et al., 2020)	1	2.92	9.83	Hat EBM (Hill et al., 2022)	51	19.30	–
DDGAN (Xiao et al., 2022)	4	3.75	9.63	CLEL-Large (Lee et al., 2023)	1200	8.61	–
ACT (Kong et al., 2024)	1	6.0	9.15	Dual-MCMC (Cui & Han, 2023)	31	9.26	8.55
<b>Diffusion-based</b>				DDAEBM (Geng et al., 2024)	4	4.82	8.86
NCSN-v2 (Song & Ermon, 2020)	1000	10.87	8.40	CDRL (Zhu et al., 2024)	96	4.31	–
DDPM (Ho et al., 2020)	1000	3.17	9.46	EC-VAE (Luo et al., 2024)	1	5.20	–
NCSN++ (Song et al., 2021)	2000	2.20	<b>9.89</b>	<b>Ours</b>			
EDM (Karras et al., 2022)	35	<b>2.04</b>	9.84	LIEBM-EM w/o MCMC	1	4.96	9.82
Flow Matching (Lipman et al., 2023)	142	6.35	–	LIEBM-EM	16	<b>4.26</b>	<b>10.02</b>
Consistency Models (Song et al., 2023)	1	8.70	8.49	LIEBM-ERM w/o MCMC	1	6.16	9.41
				LIEBM-ERM	16	4.96	9.64

## 4 EXPERIMENTS

We conduct comprehensive experiments to evaluate our proposed method under various scenarios, including unconditional image generation, OOD detection, conditional sampling, and zero-shot image restoration. We consider two options for  $\mathcal{V}$  for sampling  $p_{\text{data}}(z|x)$ : (i) the standard random augmentations commonly used in self-supervised representation learning, and (ii) adding minor uniform noise via  $x = \frac{255}{256}x + z$ , where  $z \sim \mathcal{U}(0, \frac{1}{256})$ . Both choices yield similar performance, and we adopt the first one to introduce more stochasticity into the latent variables. For our pretrained latent encoder, we evaluate three normalized self-supervised representation learning methods: SimCLR (Chen et al., 2020), BYOL (Grill et al., 2020), and W-MSE (Ermolov et al., 2021). We select SimCLR<sup>1</sup> since it achieves the best generation performance on CIFAR-10. We adopt the same architecture as Dual-MCMC for our generator and inference model. We use the energy function backbone from Dual-MCMC as our  $f_\phi$  in  $E_\theta(x, z)$ , while our  $g_\psi$  in  $p_{\phi, \psi}(z|x)$  follows the projection head architecture of the latent encoder, with Batch Normalization (Ioffe & Szegedy, 2015) removed. We apply Exponential Moving Average (EMA) with a decay rate of 0.9999 to improve generation

<sup>1</sup>We implement SimCLR using the official code of W-MSE: <https://github.com/htdt/self-supervised>

quality. We denote our model with EM generator training as LIEBM-EM, and ERM as LIEBM-ERM. For the EM setting, MCMC with marginal energy (Eq.9) outperforms joint energy matching (Eq.10), so we use Eq.9 for LIEBM-EM.

Table 2: Generative performance on CelebA 64

Model	FID ↓
SN-GAN (Miyato et al., 2018)	6.1
COCO-GAN (Lin et al., 2019)	4.0
NVAE (Vahdat & Kautz, 2020)	14.74
NCSNv2 (Song & Ermon, 2020)	26.86
DDPM (Ho et al., 2020)	3.93
<b>EBM-based</b>	
DRL (Gao et al., 2021)	5.98
VAEBM (Xiao et al., 2021)	5.31
CLEL (Lee et al., 2023)	8.34
Dual MCMC (Cui & Han, 2023)	5.15
EC-VAE (Luo et al., 2024)	<b>2.71</b>
LIEBM-EM	3.44
LIEBM-ERM	2.97

Table 3: Generative performance on CelebA-HQ 256.

Model	FID ↓
GLOW (Kingma & Dhariwal, 2018)	68.93
NVAE (Vahdat & Kautz, 2020)	45.11
VQGAN (Esser et al., 2021)	10.2
DDGAN (Xiao et al., 2022)	7.64
Score SDE (Song et al., 2021)	<b>7.23</b>
<b>EBM-based</b>	
VAEBM (Xiao et al., 2021)	20.38
Dual MCMC (Cui & Han, 2023)	15.89
CDRL (Zhu et al., 2024)	10.74
EC-VAE (Luo et al., 2024)	12.35
LIEBM-EM	10.08
LIEBM-ERM	8.76

Table 4: Generative performance on ImageNet 32.

Model	FID ↓
PixelCNN (Oord et al., 2016)	40.51
DDPM++ (Kim et al., 2021)	8.42
Flow Matching (Lipman et al., 2023)	5.02
<b>EBM-based</b>	
CF-EBM (Zhao et al., 2020)	26.31
EBM-CD (Du et al., 2021)	32.48
CLEL-Large (Lee et al., 2023)	15.47
CDRL (Zhu et al., 2024)	9.35
EC-VAE (Luo et al., 2024)	5.76
EBM <sub>MI+diff</sub> (Geng et al., 2025)	6.57
LIEBM-EM	<b>4.54</b>
LIEBM-ERM	4.98

#### 4.1 UNCONDITIONAL IMAGE GENERATION

We showcase our model’s capabilities in unconditional image generation on standard datasets involving CIFAR-10 (Krizhevsky et al., 2009), ImageNet 32 (Deng et al., 2009), CelebA 64 (Liu et al., 2015b), and CelebA-HQ 256 (Liu et al., 2015a). For quantitative results, we adopt the commonly used Fréchet inception distance (FID) and Inception Score (IS) to evaluate sample fidelity and the number of function evaluations (NFE) to evaluate sampling efficiency. We show qualitative results in Fig.3 and quantitative results in Tabs.1-4<sup>2</sup>. Fig.2 shows FID vs. network scale on CIFAR-10. Fig.5 further illustrates our method’s sampling efficiency by comparing inference time and FID across different generative models.

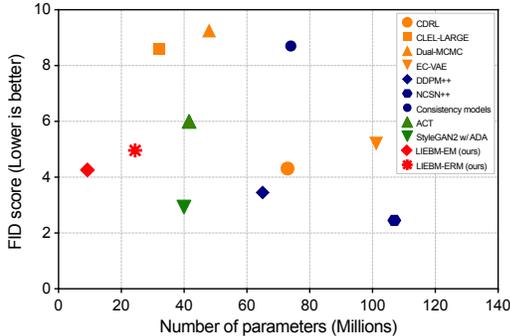


Figure 2: Param count vs. FID on CIFAR-10.

Our model achieves optimal results on most datasets, with near-optimal performance on CelebA 64 among EBMs. On CIFAR-10, our LIEBM-EM outperforms state-of-the-art CDRL, despite CDRL having 5× the number of parameters and requiring 6× MCMC steps.

Our method achieves significant improvements over CLEL with much faster sampling, demonstrating the effectiveness of our carefully designed collaborative training between the energy function and generator. Our model also outperforms Dual-MCMC and EC-VAE by a large margin on most datasets with fewer network parameters and MCMC steps, validating that our latent-informed scheme can further improve generation. Notably, for single-step generation, our model surpasses strong diffusion baselines, including Consistency Model and its adversarial variant ACT. Moreover, our model achieves competitive performance with advanced GANs and Diffusion Models while using 5-10× fewer parameters. Our model gets the best IS score on CIFAR-10 and is the first EBM to beat Flow Matching on ImageNet 32.

<sup>2</sup>Since baselines for ImageNet 32, CelebA 64, and CelebA-HQ 256 are less established than CIFAR-10, we compare using FID and commonly reported baselines.

Table 5: Inference time vs. FID on CIFAR-10.

Method	Time(s)↓	FID↓	GPU-Type
NVAE	0.36	50.97	V100
StyleGAN2 w/ ADA	0.04	2.92	V100
DDPM	80.5	3.17	V100
NCSN++	423.2	2.20	V100
EBM-Based			
VAEBM	8.79	12.19	V100
EC-VAE	0.21	5.20	RTX 2080 Ti
CLEL	82.05	8.61	RTX 2080 Ti
Dual MCMC	9.32	9.26	RTX 2080 Ti
LIEBM w/o MCMC(ours)	0.08	4.96	RTX 2080 Ti
LIEBM(ours)	1.24	4.26	RTX 2080 Ti

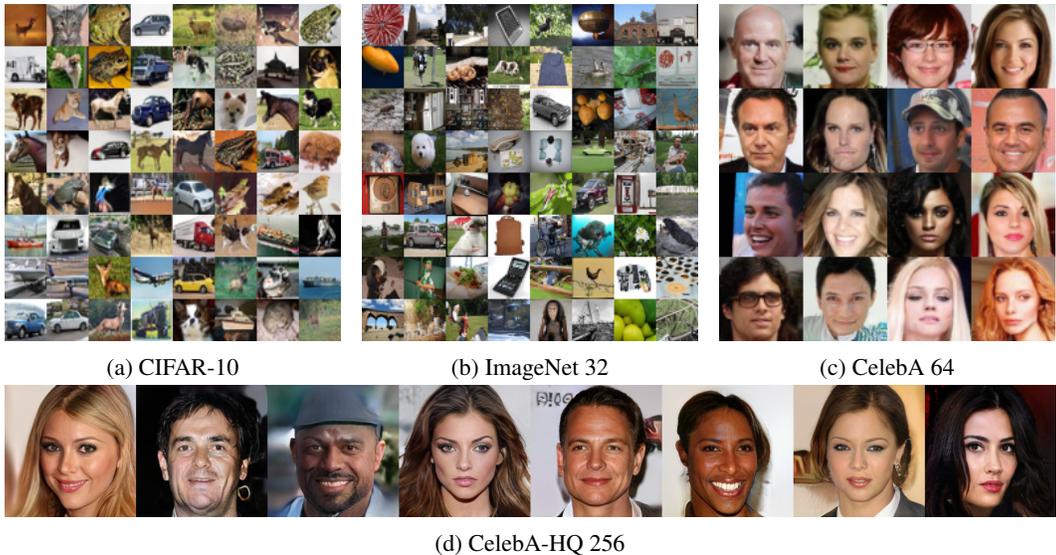


Figure 3: Samples generated by LIEBM with MCMC refinement. Select models based on FID: LIEBM-EM for CIFAR-10/ImageNet-32; LIEBM-ERM for CelebA-64/CelebA-HQ-256.

#### 4.2 OUT-OF-DISTRIBUTION DETECTION

We evaluate our model’s density modeling through OOD detection on CIFAR-10 and ImageNet 32, using their unseen test sets as inliers and other datasets as outliers. We use the standard AUROC metric with a joint energy score inspired by CLEL:

$$s(x) := F(f_{\theta}(x)) + \gamma \text{sim}(g_{\psi}(f_{\phi}(x)), h(x)). \tag{15}$$

Results are shown in Tabs.6 and 7. We observe that the joint energy score improves OOD detection on most datasets (with only slight degradation on SVHN), demonstrating enhanced robustness to diverse OOD samples through joint space modeling. On CIFAR-10, our model consistently performs at the top tier among EBMs and matches specialized OOD methods. Notably, our model shows significant improvement on CIFAR-100, which is challenging due to the similarity between CIFAR-100 and CIFAR-10. We reproduce Dual-MCMC and hat-EBM using their energy outputs as AUROC decision values. Our model exhibits strong performance on the challenging SVHN and Constant datasets for ImageNet 32, where likelihood-based methods such as VAE, GLOW, and PixelCNN typically fail at outlier detection.

Table 6: AUROC with CIFAR-10 as in-distribution.

Method	SVHN	Constant	CIFAR-100	CelebA
PixelCNN++ (Salimans et al., 2017)	0.32	0.71	0.63	-
GLOW (Kingma & Dhariwal, 2018)	0.24	-	0.55	0.57
NVAE (Vahdat & Kautz, 2020)	0.44	0.65	0.49	0.68
JEM (Duvenaud et al., 2020)	0.67	-	0.67	0.75
DRL (Gao et al., 2021)	0.88	0.99	0.44	0.64
hatEBM (Hill et al., 2022)	0.75	0.36	0.63	0.62
CLEL (Lee et al., 2023)	0.98	-	0.72	0.77
Dual-MCMC(Cui & Han, 2023)	0.62	0.32	0.54	0.59
<b>Specialized OOD methods</b>				
OOD EBM (Liu et al., 2020)	0.91	-	<b>0.87</b>	<b>0.78</b>
MPDR-S (Yoon et al., 2023)	<b>0.99</b>	<b>0.9996</b>	0.56	0.73
LIEBM-EM ( $f_\theta(x)$ )	0.96	0.67	0.66	0.68
LIEBM-EM	0.95	0.97	0.82	0.77
LIEBM-ERM ( $f_\theta(x)$ )	0.94	0.76	0.68	0.58
LIEBM-ERM	0.95	0.96	0.82	0.75

Table 7: AUROC with ImageNet 32 as in-distribution. ( $f_\theta(x)$ ) means  $f_\theta(x)$  serves as the decision function.

Method	SVHN	Constant	FMNIST	CelebA
DAE (Vincent et al., 2008)	0.10	0.07	0.991	0.43
VAE (Kingma & Welling, 2014)	0.13	0.03	0.95	0.55
WAE (Tolstikhin et al., 2018)	0.08	0.07	0.991	0.36
PixelCNN++ (Salimans et al., 2017)	0.03	0.00	0.004	0.24
GLOW (Kingma & Dhariwal, 2018)	0.17	0.41	0.86	0.48
CLEL (Lee et al., 2023)	0.96	0.83	0.54	0.74
<b>Specialized OOD methods</b>				
NAE (Yoon et al., 2021)	0.985	0.97	<b>0.994</b>	<b>0.95</b>
LIEBM-EM ( $f_\theta(x)$ )	<b>0.99</b>	0.97	0.40	0.48
LIEBM-EM	0.984	<b>0.99</b>	0.896	0.52
LIEBM-ERM ( $f_\theta(x)$ )	<b>0.99</b>	0.93	0.35	0.45
LIEBM-ERM	0.985	<b>0.99</b>	0.868	0.54

### 4.3 CONDITIONAL SAMPLING

We also investigate conditional sampling with our latent representation as labels. Unlike CLEL, we employ a generator as an initializer, which offers faster sampling but requires the generator to produce high-quality initial samples. Therefore, similar to the ERM setting, we train an inference model to form an autoencoder with the generator under our EM framework, enabling us to obtain reconstructions from the input for initialization. We train our inference model using a variant of ELBO loss in the latent space to ensure detailed clarity and sharpness while preserving semantic similarity (See Appendix A.6 for more results). Specifically, we use Eq.14 to train our inference model, but omitting the pixel-level reconstruction term  $\log p_g(x|m)$ . We split our generation into two components  $G(m) + Y$ . Following CLEL, we obtain the class representation  $\bar{z}_c$  for each class  $c$ , defined as the normalized average of latent representation across all images in class  $c$ . We draw an initialization  $\bar{x}_c$  as initial  $G(m)$  by averaging all augmented images from each class. Then we iteratively optimize  $Y$  and  $m$  by performing MCMC sampling from  $E_\theta(G(m) + Y, \bar{z}_c)$  and using the inference model conditioned on  $G(m) + Y$ , respectively. From Fig.4 we can see that the EM setting is able to generate diverse samples with clear details for each class, whereas the ERM setting, while capable of generating some feature elements of the given class, fails to produce identifiable subjects. This is caused by the ELBO component in ERM training, which provides pixel-level reconstruction but produces blurry, low-sharpness results.

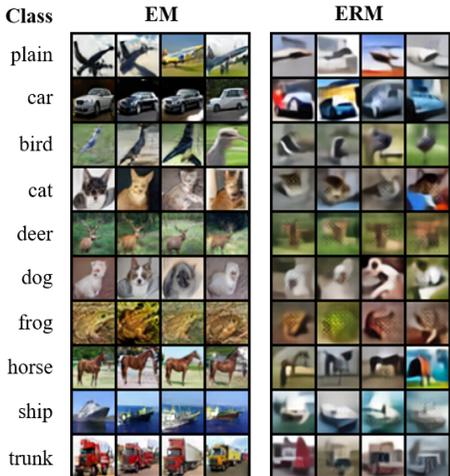


Figure 4: Conditional generated sample on CIFAR-10.

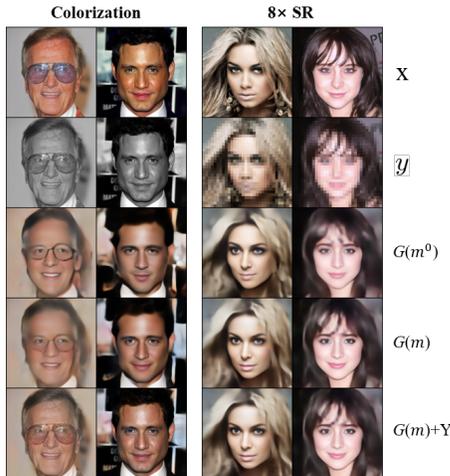


Figure 5: Qualitative results of zero-shot image restoration on CelebA-HQ 256.

#### 4.4 IMAGE RESTORATION

We also present the application of our method in zero-shot image restoration tasks, including colorization and  $8\times$  super-resolution. We conduct experiments on CelebA-HQ 256 with ERM setting, since we need pixel-level restoration. Following Luo et al. (2024); Wang et al. (2023), we also use a linear operator  $A$  to get the degraded image  $y = Ax$  and utilize its pseudo-inverse  $A^\dagger$  to derive the initial estimate  $\hat{x} = A^\dagger y$ . We obtain initial  $m^0$  using our inference model with input  $\hat{x}$ . Inspired by Luo et al. (2024), we use the following joint function to refine  $m$  using MCMC sampling:

$$p_{g,\theta}(A^\dagger y, m) \propto \exp(E_\theta(G(m), h(G(m)))) p(m) p(A^\dagger y | A^\dagger A G(m)) \quad (16)$$

After refining  $m$ , we update  $Y$  using MCMC sampling with  $G(m)$  fixed:

$$p_{g,\theta}(A^\dagger y, Y) \propto \exp(E_\theta(G(m) + Y, h(G(m) + Y))) p(A^\dagger y | A^\dagger A (G(m) + Y)) \quad (17)$$

We employ  $\tilde{x} = G(m) + Y$  as our restoration solution. The qualitative results are shown in Fig.5 and the corresponding PSNR and SSIM metrics are reported in Tab.8. We can observe that with the help of joint energy distribution, our model can successfully restore those images with high quality and consistency after refinement on  $m$  and  $Y$ .

#### 4.5 ABLATION STUDY

Fig.6 tracks FID scores during training for the ablation study on CIFAR-10. Tab.9 shows their corresponding OOD performance. It can be seen that traditional EBM training without latent variables can not converge, no matter how the generator training is designed. Training with a pretrained latent encoder improves both generation performance and OOD robustness while yielding a better latent encoder with enhanced semantic separability, as shown in Fig.7. Augmentation technique can improve OOD results on Constant Dataset with negligible generation degradation. Generator training with Eq.10 for MCMC refinement (EJM) can get better results for the first stage, but finally slightly worse than EM setting (Eq.9). In particular, EJM tends to collapse towards the end of training on ImageNet 32. Hence, we recommend employing the EM setting. Combining Tabs.1-4, we can see that our EM setting achieves better performance than ERM on multi-class datasets such as CIFAR-10 and ImageNet, while for few-modal datasets like CelebA and CelebA-HQ, ERM performs better.

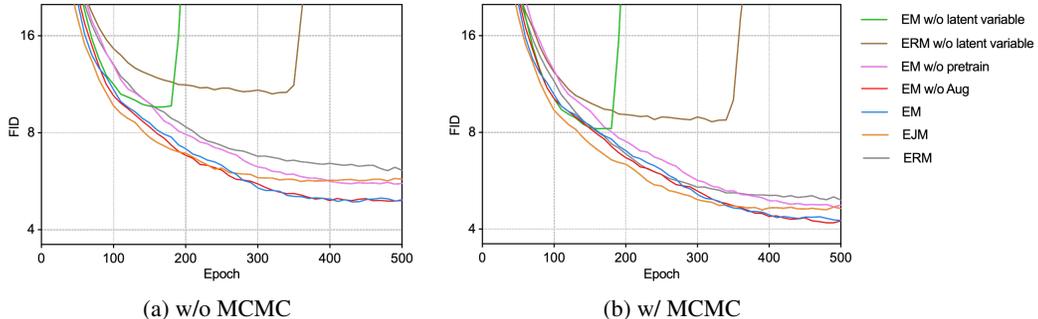


Figure 6: FID with different settings on CIFAR-10. “w/o MCMC” means direct sampling from the generator without MCMC refinement.

**Generality with different  $\mathcal{V}$  choices and energy forms** Fig.8 compares different  $\mathcal{V}$  choices and energy forms. For  $\mathcal{V}$  choices, we observe similar performance regardless of using random augmentations or uniform noise. This shows that our method is agnostic to the design of  $\mathcal{V}$ . For energy forms, while the Gaussian distribution is the conventional choice for modeling explicit posterior

Table 8: Quantitative results of zero-shot image restoration on CelebA-HQ 256.

Model	Colorization		$8\times$ SR	
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$
$G(m^0)$	20.64	0.66	22.62	0.67
$G(m)$	22.02	0.70	24.16	0.70
$G(m) + Y$	25.25	0.94	24.55	0.71

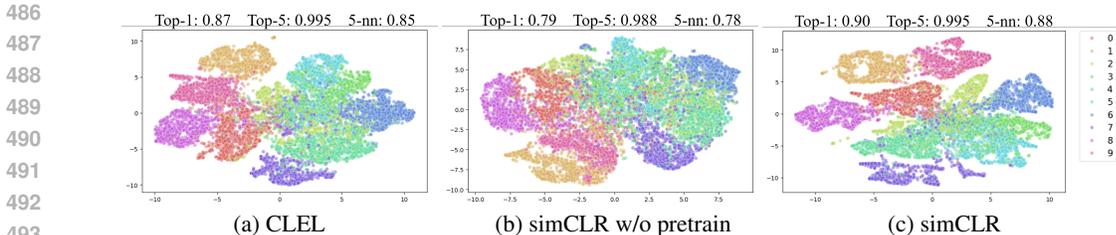


Figure 7: t-SNE visualization of latent representation on CIFAR-10 test set. Accuracies of a linear classifier (Top-1 & Top-5) and a 5-nearest neighbors classifier are shown above each subfigure.

Table 9: AUROC under different settings with CIFAR-10 as in-distribution.

Method	SVHN	Constant	CIFAR-100	CelebA
EM w/o pretrain	<b>0.97</b>	<b>0.997</b>	0.75	0.71
EM w/o Aug	0.95	0.88	<b>0.82</b>	0.76
EM	0.95	0.97	<b>0.82</b>	<b>0.77</b>
EJM	0.95	0.94	0.81	0.73
ERM	0.95	0.96	<b>0.82</b>	0.75

Table 10: Performance with different normalized SSRL methods.

Method	FID	AUROC			
		SVHN	Constant	FMNIST	CelebA
BYOL	5.23	<b>0.96</b>	0.98	<b>0.85</b>	<b>0.81</b>
W-MSE	5.16	0.93	<b>0.99</b>	0.83	0.77
SimCLR	<b>4.26</b>	0.95	0.97	0.82	0.77

$p_{\phi, \psi}(z|x)$  (Han et al., 2019; Kan et al., 2022), we observe its sensitive variance effects. Thus, we constrain the normalized mean and log-variance to  $[-1, 1]$  for the Gaussian posterior, and fix the variance for CEBM. Both perform better than the Gaussian-posterior Dual MCMC, yet remain inferior to the cosine-similarity posterior. Moreover, by fixing the variance, the squared Euclidean distance  $\|z_1 - z_2\|_2^2 = 2 - 2 \text{sim}(z_1, z_2)$ , cosine-similarity form naturally induces a spherical Gaussian on the unit sphere. This shows that using a cosine-similarity posterior and decoupling the implicit data energy from the explicit posterior, as in Eq.5, offers greater flexibility and easier control.

**Adaptability to various self-supervised representation learning methods.** Our framework theoretically can be applied to any normalized self-supervised representation learning (SSRL) method. To verify our model’s adaptability, we choose two other classic normalized SSRL methods, BYOL and W-MSE, to pretrain our latent encoder. Tab.10 reports FID and AUROC metrics for different SSRL methods, confirming that our LIEBM scales well to various SSRL methods.

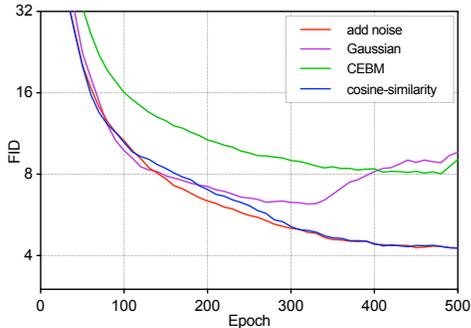


Figure 8: FID with different  $\mathcal{V}$  choices and energy forms on CIFAR-10. “add noise” means adding uniform noise to construct  $\mathcal{V}$ . “Gaussian” and “cosine-similarity” mean two posterior choices.

## 5 CONCLUSION

In this paper, we propose LIEBM, a collaborative training scheme that jointly learns a latent-informed EBM and its generator initializer. We leverage pretrained self-supervised representations as our target latent variables to guide the energy function in capturing the semantic structure of the data manifold. Our model narrows the gap between EBMs and mainstream generative models while retaining the benefits of lightweight architectures. It also excels in various downstream tasks, such as OOD detection, conditional sampling, and zero-shot image restoration. Additionally, our framework could be extended to multi-modal large models by treating the joint space as a multi-modal space and replacing SSRL methods with advanced modal-alignment techniques such as CLIP and ALBEF. We hope our work brings to light the profound potential of EBMs as mainstream generative models and stimulate active research in this area.

## REFERENCES

- 540  
541  
542 Michal Balcerak, Tamaz Amiranashvili, Antonio Terpin, Suprosanna Shit, Lea Bogensperger, Se-  
543 bastian Kaltenbach, Petros Koumoutsakos, and Bjoern Menze. Energy matching: Unifying flow  
544 matching and energy-based models for generative modeling. *Advances in Neural Information  
545 Processing Systems*, 2025.
- 546 Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural  
547 image synthesis. *International Conference on Learning Representations*, 2019.
- 548  
549 Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of markov chain monte  
550 carlo*. CRC press, 2011.
- 551  
552 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
553 contrastive learning of visual representations. *International conference on machine learning*, pp.  
554 1597–1607, 2020.
- 555  
556 Jiali Cui and Tian Han. Learning energy-based model via dual-mcmc teaching. *Advances in Neural  
557 Information Processing Systems*, 36:28861–28872, 2023.
- 558  
559 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-  
560 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
pp. 248–255. IEEE, 2009. doi: 10.1109/CVPR.2009.5206848.
- 561  
562 Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *Advances  
563 in Neural Information Processing Systems*, 32, 2019.
- 564  
565 Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based mod-  
566 els. *Advances in Neural Information Processing Systems*, 33:6637–6647, 2020.
- 567  
568 Yilun Du, Shuang Li, Joshua Tenenbaum, and Igor Mordatch. Improved contrastive divergence  
569 training of energy based models. *International Conference on Machine Learning*, 2021.
- 570  
571 Yilun Du, Conor Durkan, Robin Strudel, Joshua B Tenenbaum, Sander Dieleman, Rob Fergus,  
572 Jascha Sohl-Dickstein, Arnaud Doucet, and Will Sussman Grathwohl. Reduce, reuse, recycle:  
573 Compositional generation with energy-based diffusion models and mcmc. *International confer-  
574 ence on machine learning*, pp. 8489–8510, 2023.
- 575  
576 Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Alex Lamb, Martin Arjovsky, Olivier Mastropi-  
577 etro, and Aaron Courville. Adversarially learned inference. *International Conference on Learning  
578 Representations*, Jun 2016.
- 579  
580 David Duvenaud, Jackson Wang, Jorn Jacobsen, Kevin Swersky, Mohammad Norouzi, and Will  
581 Grathwohl. Your classifier is secretly an energy based model and you should treat it like one.  
582 *International Conference on Learning Representations*, 4, 2020.
- 583  
584 Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-  
585 supervised representation learning. *International Conference on Machine Learning*, pp. 3015–  
586 3024, 2021.
- 587  
588 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image  
589 synthesis. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
590 Jun 2021. doi: 10.1109/cvpr46437.2021.01268. URL [http://dx.doi.org/10.1109/  
591 cvpr46437.2021.01268](http://dx.doi.org/10.1109/cvpr46437.2021.01268).
- 592  
593 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*, pp. 7518–7528, 2020.
- Ruiqi Gao, Yang Song, Ben Poole, Ying Nian Wu, and Diederik P Kingma. Learning energy-based  
models by diffusion recovery likelihood. *International Conference on Learning Representations*,  
2021.

- 594 Cong Geng, Jia Wang, Zhiyong Gao, Jes Frellsen, and Søren Hauberg. Bounds all around: training  
595 energy-based models with bidirectional bounds. *Advances in Neural Information Processing*  
596 *Systems*, 34:19808–19821, 2021.
- 597 Cong Geng, Jia Wang, Li Chen, and Zhiyong Gao. Solving the reconstruction-generation trade-off:  
598 Generative model with implicit embedding learning. *Neurocomputing*, 549:126428, 2023.
- 600 Cong Geng, Tian Han, Peng-Tao Jiang, Hao Zhang, Jinwei Chen, Søren Hauberg, and Bo Li. Im-  
601 proving adversarial energy-based model via diffusion process. *International Conference on Ma-*  
602 *chine Learning*, 2024.
- 603 Cong Geng, Jia Wang, Zhiyong Gao, Jes Frellsen, and Søren Hauberg. Exploring bidirectional  
604 bounds for minimax-training of energy-based models. *International Journal of Computer Vision*,  
605 May 2025. doi: 10.1007/s11263-025-02460-0.
- 607 Will Sussman Grathwohl, Jacob Jin Kelly, Milad Hashemi, Mohammad Norouzi, Kevin Swersky,  
608 and David Duvenaud. No MCMC for me: Amortized sampling for fast and stable training of  
609 energy-based models. *International Conference on Learning Representations*, 2021.
- 610 Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
611 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
612 et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in Neural*  
613 *Information Processing Systems*, 33:21271–21284, 2020.
- 614 Tian Han, Erik Nijkamp, Xiaolin Fang, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Diver-  
615 gence triangle for joint training of generator model, energy-based model, and inferential model.  
616 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
617 8670–8679, 2019.
- 619 Tian Han, Erik Nijkamp, Linqi Zhou, Bo Pang, Song-Chun Zhu, and Ying Nian Wu. Joint training  
620 of variational auto-encoder and latent energy-based model. In *Proceedings of the IEEE/CVF*  
621 *conference on computer vision and pattern recognition*, pp. 7978–7987, 2020.
- 622 Mitch Hill, Erik Nijkamp, Jonathan Mitchell, Bo Pang, and Song-Chun Zhu. Learning probabilistic  
623 models from generator latent spaces with hat ebm. *Advances in Neural Information Processing*  
624 *Systems*, 35:928–940, 2022.
- 625 Geoffrey E. Hinton, Peter Dayan, Brendan J. Frey, and Radford M. Neal. The “wake-sleep” al-  
626 gorithm for unsupervised neural networks. *Science*, pp. 1158–1161, May 1995. doi: 10.1126/  
627 science.7761831. URL <http://dx.doi.org/10.1126/science.7761831>.
- 629 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
630 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 631 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by  
632 reducing internal covariate shift. *International Conference on Machine Learning*, pp. 448–456,  
633 2015.
- 634 Ge Kan, Jinhu Lü, Tian Wang, Baochang Zhang, Aichun Zhu, Lei Huang, Guodong Guo, and  
635 Hichem Snoussi. Bi-level doubly variational learning for energy-based latent variable models.  
636 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.  
637 18460–18469, 2022.
- 639 Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training  
640 generative adversarial networks with limited data. *Advances in Neural Information Processing*  
641 *Systems*, 33:12104–12114, 2020.
- 642 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-  
643 based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577,  
644 2022.
- 645 Dongjun Kim, Seungjae Shin, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Soft truncation:  
646 A universal training technique of score-based diffusion model for high precision score estimation.  
647 *International Conference on Machine Learning*, Jun 2021.

- 648 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *International Conference*  
649 *on Learning Representations*, 2014.
- 650
- 651 Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions.  
652 *Advances in Neural Information Processing Systems*, 31, 2018.
- 653 Fei Kong, Jinhao Duan, Lichao Sun, Hao Cheng, Renjing Xu, Hengtao Shen, Xiaofeng Zhu, Xi-  
654 aoshuang Shi, and Kaidi Xu. Act-diffusion: Efficient adversarial consistency training for one-step  
655 diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
656 *Recognition*, pp. 8890–8899, 2024.
- 657
- 658 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
659 2009.
- 660 Rithesh Kumar, Sherjil Ozair, Anirudh Goyal, Aaron Courville, and Yoshua Bengio. Maximum  
661 entropy generators for energy-based models. *arXiv preprint arXiv:1901.08508*, 2019.
- 662
- 663 Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, Fugie Huang, et al. A tutorial on energy-  
664 based learning. *Predicting structured data*, 1(0), 2006.
- 665 Hankook Lee, Jongheon Jeong, Sejun Park, and Jinwoo Shin. Guiding energy-based models via  
666 contrastive latent variables. *International Conference on Learning Representations*, 2023.
- 667
- 668 Chunyuan Li, Hao Liu, Changyou Chen, Yunchen Pu, Liqun Chen, Ricardo Henao, and Lawrence  
669 Carin. Alice: Towards understanding adversarial learning for joint distribution matching. *Neural*  
670 *Information Processing Systems*, Sep 2017.
- 671 Chieh Hubert Lin, Chia-Che Chang, Yu-Sheng Chen, Da-Cheng Juan, Wei Wei, and Hwann-Tzong  
672 Chen. Coco-gan: Generation by parts via conditional coordinating. In *Proceedings of the*  
673 *IEEE/CVF international conference on computer vision*, pp. 4512–4521, 2019.
- 674
- 675 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
676 for generative modeling. *International Conference on Learning Representations*, 2023.
- 677 Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detec-  
678 tion. *Advances in Neural Information Processing Systems*, 33:21464–21475, 2020.
- 679
- 680 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.  
681 In *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec 2015a. doi: 10.1109/  
682 iccv.2015.425. URL <http://dx.doi.org/10.1109/iccv.2015.425>.
- 683 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.  
684 In *International Conference on Computer Vision*, pp. 3730–3738, 2015b.
- 685
- 686 Junn Yong Loo, Leong Fang Yu, Michelle Adeline, Julia K Lau, Hwa Hui Tew, Arghya Pal,  
687 Vishnu Monn Baskaran, Chee-Ming Ting, and Raphael CW Phan. Learning energy-based gener-  
688 ative models via potential flow: A variational principle approach to probability density homotopy  
689 matching. *Transactions on Machine Learning Research*, 2025.
- 690 Yihong Luo, Siya Qiu, Xingjian Tao, Yujun Cai, and Jing Tang. Energy-calibrated vae with test time  
691 free lunch. In *European Conference on Computer Vision*, pp. 326–344. Springer, 2024.
- 692
- 693 Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for  
694 generative adversarial networks. *International Conference on Learning Representations*, 2018.
- 695
- 696 Radford M Neal et al. Mcmc using hamiltonian dynamics. *Handbook of markov chain monte carlo*,  
2(11):2, 2011.
- 697
- 698 Kirill Neklyudov, Rob Brekelmans, Daniel Severo, and Alireza Makhzani. Action matching: Learn-  
699 ing stochastic dynamics from samples. In *International conference on machine learning*, pp.  
25858–25889. PMLR, 2023.
- 700
- 701 Erik Nijkamp, Mitch Hill, Song-Chun Zhu, and Ying Wu. Learning non-convergent non-persistent  
short-run mcmc toward energy-based model. *Neural Information Processing Systems*, Jan 2019.

- 702 Erik Nijkamp, Mitch Hill, Tian Han, Song-Chun Zhu, and Ying Nian Wu. On the anatomy of  
703 mcmc-based maximum likelihood learning of energy-based models. In *Proceedings of the AAAI*  
704 *Conference on Artificial Intelligence*, volume 34, pp. 5272–5280, 2020.
- 705  
706 Erik Nijkamp, Ruiqi Gao, Pavel Soutsov, Srinivas Vasudevan, Bo Pang, Song-Chun Zhu, and  
707 Ying Nian Wu. Mcmc should mix: Learning energy-based model with neural transport latent  
708 space mcmc. *International Conference on Learning Representations*, 2022.
- 709  
710 Aaronvanden Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks.  
711 *International Conference on Machine Learning*, Jan 2016.
- 712  
713 Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic  
714 gradient langevin dynamics: a nonasymptotic analysis. In *Conference on Learning Theory*, pp.  
715 1674–1703. PMLR, 2017.
- 716  
717 Ruslan Salakhutdinov, Andriy Mnih, and Geoffrey Hinton. Restricted boltzmann machines for col-  
718 laborative filtering. In *Proceedings of the 24th international conference on Machine learning*, Jun  
719 2007. doi: 10.1145/1273496.1273596. URL <http://dx.doi.org/10.1145/1273496.1273596>.
- 720  
721 Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the  
722 pixelcnn with discretized logistic mixture likelihood and other modifications. *International Con-*  
723 *ference on Learning Representations*, 2017.
- 724  
725 Yang Song and Stefano Ermon. Improved techniques for training score-based generative models.  
726 *Advances in Neural Information Processing Systems*, 33:12438–12448, 2020.
- 727  
728 Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint*  
729 *arXiv:2101.03288*, 2021.
- 730  
731 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
732 Poole. Score-based generative modeling through stochastic differential equations. *International*  
733 *Conference on Learning Representations*, 2021.
- 734  
735 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *International*  
736 *Conference on Machine Learning*, 2023.
- 737  
738 Yunfu Song and Zhijian Ou. Generative modeling by inclusive neural random fields with applica-  
739 tions in image generation and anomaly detection. *arXiv preprint arXiv:1806.00271*, 2018.
- 740  
741 Tijmen Tieleman. Training restricted boltzmann machines using approximations to the likelihood  
742 gradient. *International Conference on Machine Learning*, Jan 2008.
- 743  
744 Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-  
745 encoders. *International Conference on Learning Representations*, 2018.
- 746  
747 Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. *Advances in Neural*  
748 *Information Processing Systems*, 33:19667–19679, 2020.
- 749  
750 Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and  
751 composing robust features with denoising autoencoders. *International Conference on Machine*  
752 *Learning*, Jan 2008. doi: 10.1145/1390156.1390294. URL <http://dx.doi.org/10.1145/1390156.1390294>.
- 753  
754 Ben Wan, Cong Geng, Tianyi Zheng, and Jia Wang. Ebm-wgf: Training energy-based models with  
755 wasserstein gradient flow. *Neural Networks*, 187:107300, 2025.
- 756  
757 Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion  
758 null-space model. *International Conference on Learning Representations*, 2023.
- 759  
760 Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. *Inter-*  
761 *national Conference on Machine Learning*, pp. 681–688, 2011.

- 756 Hao Wu, Babak Esmaeili, Michael Wick, Jean-Baptiste Tristan, and Jan-Willem Van De Meent.  
757 Conjugate energy-based models. *International Conference on Machine Learning*, pp. 11228–  
758 11239, 2021.
- 759
- 760 Zhisheng Xiao, Karsten Kreis, Jan Kautz, and Arash Vahdat. Vaebm: A symbiosis between varia-  
761 tional autoencoders and energy-based models. *International Conference on Learning Representations*, 2021.  
762
- 763
- 764 Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with  
765 denoising diffusion gans. *International Conference on Learning Representations*, 2022.
- 766
- 767 Jianwen Xie, Song-Chun Zhu, and Ying Nian Wu. Learning energy-based spatial-temporal genera-  
768 tive convnets for dynamic patterns. *IEEE Transactions on Pattern Analysis and Machine Intelli-*  
769 *gence*, 43(2):516–531, 2019.
- 770
- 771 Jianwen Xie, Yang Lu, Ruiqi Gao, Song-Chun Zhu, and Ying Nian Wu. Cooperative training of  
772 descriptor and generator networks. *IEEE Transactions on Pattern Analysis and Machine Intelli-*  
773 *gence*, 42(1):27–45, 2020.
- 774
- 775 Jianwen Xie, Zilong Zheng, Xiaolin Fang, Song-Chun Zhu, and Ying Nian Wu. Cooperative training  
776 of fast thinking initializer and slow thinking solver for conditional learning. *IEEE Transactions*  
777 *on Pattern Analysis and Machine Intelligence*, 44(8):3957–3973, 2021a.
- 778
- 779 Jianwen Xie, Zilong Zheng, and Ping Li. Learning energy-based model with variational auto-  
780 encoder as amortized sampler. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
781 volume 35, pp. 10441–10451, 2021b.
- 782
- 783 Jianwen Xie, Yaxuan Zhu, Jun Li, and Ping Li. A tale of two flows: Cooperative learning of  
784 langevin flow and normalizing flow toward energy-based model. *International Conference on*  
785 *Learning Representations*, 2022.
- 786
- 787 Pan Xu, Jinghui Chen, Difan Zou, and Quanquan Gu. Global convergence of langevin dynam-  
788 ics based algorithms for nonconvex optimization. *Advances in Neural Information Processing*  
789 *Systems*, 31, 2018.
- 790
- 791 Sangwoong Yoon, Yung-Kyun Noh, and Frank Park. Autoencoding under normalization constraints.  
792 *International Conference on Machine Learning*, pp. 12087–12097, 2021.
- 793
- 794 Sangwoong Yoon, Young-Uk Jin, Yung-Kyun Noh, and Frank Park. Energy-based models for  
795 anomaly detection: A manifold diffusion recovery approach. *Advances in Neural Information*  
796 *Processing Systems*, 36:49445–49466, 2023.
- 797
- 798 Yang Zhao, Jianwen Xie, and Ping Li. Learning energy-based generative models via coarse-to-fine  
799 expanding and sampling. *International Conference on Learning Representations*, 2020.
- 800
- 801

## 802 A APPENDIX

### 803 A.1 LLM USAGE

804 We employ large language models (LLMs) to assist with language polishing and grammar improve-  
805 ment throughout the paper. For the Related Work section, we leverage LLMs to help synthesize  
806 brief summaries of related research publications. We also use LLMs to help generate some simple  
807 experiment code and LaTeX formatting code for figures and tables. We have verified and validated  
808 all contents made by LLMs and take full responsibility for our submission.  
809

## 810 A.2 RELATED WORK

811  
812 Energy-based models (EBMs) represent a powerful class of generative models that offer explicit  
813 unnormalized density estimation and architectural flexibility. Traditional EBM training relies on  
814 maximum likelihood estimation (MLE) with Markov Chain Monte Carlo (MCMC) sampling, par-  
815 ticularly Langevin dynamics. However, noise-initialized Langevin dynamics often suffer from slow  
816 convergence and computational inefficiency (Song & Kingma, 2021). Several techniques have been  
817 proposed to alleviate the expensive MCMC, such as Persistent Contrastive Divergence (PCD) (Tiele-  
818 man, 2008), adding a replay buffer (Du & Mordatch, 2019), short-run MCMC (Nijkamp et al.,  
819 2019), et.al. Nevertheless, these approaches remain inefficient as they still require hundreds to thou-  
820 sands of MCMC steps. Cooperative learning methods (Xie et al., 2020; 2021b; 2022; Hill et al.,  
821 2022) introduce a generator as a fast initializer learned to amortize long-run MCMC. Adversarial  
822 EBMs (Kumar et al., 2019; Geng et al., 2021; Grathwohl et al., 2021; Wan et al., 2025) form a mini-  
823 max game between the energy function and the introduced generator to enable MCMC-free training.  
824 Some advances link connections between EBMs and other generative models to benefit from their  
825 strengths, such as VAE (Xiao et al., 2021; Luo et al., 2024), flow-based models (Nijkamp et al.,  
826 2022; Gao et al., 2020), and diffusion-based models (Gao et al., 2021; Zhu et al., 2024; Geng et al.,  
827 2024). Some recent works (Neklyudov et al., 2023; Loo et al., 2025; Balcerak et al., 2025) formulate  
828 EBMs from a vector field perspective, but slow sampling remains an issue.

829 Latent-variable EBMs define an energy function to characterize the joint density over data and latent  
830 variables. CLEL (Lee et al., 2023) leverages contrastive representation learning to learn meaning-  
831 ful latent structures that subsequently guide the EBM training. CEBM (Wu et al., 2021) decom-  
832 poses the joint density into an intractable data distribution and a tractable latent posterior, providing  
833 VAE-like functionality while preserving EBM interpretability and density estimation. Divergence  
834 Triangle (Han et al., 2019; 2020) and Dual-MCMC (Cui & Han, 2023) build a unified framework  
835 that employs divergence triangle formulations to seamlessly integrate energy function, generator,  
836 and inference model through minimizing KL divergences between joint distributions. We focus on  
837 collaborative learning between the generator and latent-variable EBM, decoupling the latent distri-  
838 bution from the generator’s prior to retain informative latent representations.

## 838 A.3 PRELIMINARY OF EBMS

839  
840 Let  $\mathcal{X}$  be the data space and  $p_{\text{data}}(\mathbf{x})$  be true data distribution. An EBM defines a probability distri-  
841 bution through an energy function  $E_{\theta} : \mathcal{X} \rightarrow \mathbb{R}$  parameterized by  $\theta$ ,

$$842 p_{\theta}(\mathbf{x}) = \frac{\exp(E_{\theta}(\mathbf{x}))}{Z_{\theta}}, \quad Z_{\theta} = \int \exp(E_{\theta}(\mathbf{x})) d\mathbf{x}, \quad (18)$$

843 where  $Z_{\theta}$  is the intractable normalizing constant. EBMs primarily rely on maximizing the log-  
844 likelihood for training such that:

$$845 L(\theta) := \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [E_{\theta}(\mathbf{x})] - \log Z_{\theta}. \quad (19)$$

846 The gradient of  $L(\theta)$  can be derived as:

$$847 \frac{\partial L}{\partial \theta} = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \left[ \frac{\partial}{\partial \theta} E_{\theta}(\mathbf{x}) \right] - \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x})} \left[ \frac{\partial}{\partial \theta} E_{\theta}(\mathbf{x}) \right]. \quad (20)$$

848 Eq.20 requires MCMC sampling from energy distribution  $p_{\theta}(\mathbf{x})$ , which can be achieved by Langevin  
849 dynamics (Welling & Teh, 2011):

$$850 \mathbf{x}^{t+1} = \mathbf{x}^t + \frac{\delta^2}{2} \nabla_{\mathbf{x}} E_{\theta}(\mathbf{x}^t) + \delta \epsilon^t, \quad (21)$$

851 where  $t$  indexes the time step,  $\delta$  is the step size, and  $\epsilon \sim \mathcal{N}(0, I)$ . For small enough  $\epsilon$  and large  
852 enough  $t$ , the distribution of  $\mathbf{x}^t$  weakly converges to the energy distribution  $p_{\theta}(\mathbf{x})$  regardless of the  
853 initial distribution of  $\mathbf{x}^0$  (Raginsky et al., 2017; Xu et al., 2018; Neal et al., 2011).

## 854 A.4 DERIVATION OF EQ.4

855 From Eq.3, we have

$$856 \frac{\partial L}{\partial \theta} = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim p_{\text{data}}(\mathbf{x}, \mathbf{z})} \left[ \frac{\partial}{\partial \theta} E_{\theta}(\mathbf{x}, \mathbf{z}) \right] - \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim p_{\theta}(\mathbf{x}, \mathbf{z})} \left[ \frac{\partial}{\partial \theta} E_{\theta}(\mathbf{x}, \mathbf{z}) \right] \quad (22)$$

Since  $E_\theta(\mathbf{x}) = \log \int \exp(E_\theta(\mathbf{x}, z)) dz$ , then  $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}, z) dz = \frac{\exp(E_\theta(\mathbf{x}))}{Z_\theta}$ , thus  $E_\theta(\mathbf{x})$  is an available energy function of marginal  $p_\theta(\mathbf{x})$ . We can obtain:

$$p_\theta(\mathbf{x}, z) = \frac{\exp(E_\theta(\mathbf{x}, z))}{Z_\theta} = \frac{\exp(E_\theta(\mathbf{x}))}{Z_\theta} p_\theta(z|\mathbf{x})$$

$$E_\theta(\mathbf{x}, z) = E_\theta(\mathbf{x}) + \log p_\theta(z|\mathbf{x}) \quad (23)$$

Substituting Eq.23 into the second term of Eq.22 yields:

$$\mathbb{E}_{(\mathbf{x}, z) \sim p_\theta(\mathbf{x}, z)} \left[ \frac{\partial}{\partial \theta} E_\theta(\mathbf{x}, z) \right] = \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[ \frac{\partial}{\partial \theta} E_\theta(\mathbf{x}) \right] + \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x}, z)} \left[ \frac{\partial}{\partial \theta} \log p_\theta(z|\mathbf{x}) \right]$$

$$= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[ \frac{\partial}{\partial \theta} E_\theta(\mathbf{x}) \right] \quad (24)$$

The second equality follows from:

$$\mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x}, z)} \left[ \frac{\partial}{\partial \theta} \log p_\theta(z|\mathbf{x}) \right] = \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})} \left[ \int \frac{\partial}{\partial \theta} p_\theta(z|\mathbf{x}) dz \right] = 0 \quad (25)$$

Plugging Eq.24 in Eq.22, we can get Eq.4.

## A.5 TRAINING PROCEDURE OF LIEBM

---

### Algorithm 1 LIEBM training

---

**Require:** a latent-variable EBM  $E_\theta(\mathbf{x}, z)$ , a Generator  $G$ , an inference model  $q_\alpha$ , a pretrained latent encoder  $h$ , an augmentation distribution  $\mathcal{V}$ , hyperparameters  $\tau, p$ .

- 1: **for** # training iterations **do**
  - 2:   Sample  $\{\mathbf{x}_i\}_{i=1}^n \sim p_{\text{data}}(\mathbf{x})$ .
  - 3:   Sample  $z_i = h(v(\mathbf{x}_i)/\|h(v(\mathbf{x}_i))\|_2)$ ,  $v \sim \mathcal{V}$  to get sample pairs  $\{\mathbf{x}_i, z_i\}_{i=1}^n \sim p_{\text{data}}(\mathbf{x}, z)$   
     ▷ Generator training
  - 4:   Sample generated examples  $\{G(\mathbf{m}_i)\}_{i=1}^n$ , where  $\mathbf{m}_i \sim \mathcal{N}(0, I)$
  - 5:   Obtain refined samples  $\{\tilde{\mathbf{x}}_i^T\}_{i=1}^n \sim p_\theta(\mathbf{x})$  initialized from  $\tilde{\mathbf{x}}_i^0 = G(\mathbf{m}_i)$  using Eq.9 or Eq.10.
  - 6:   Compute  $L_G = \frac{1}{n} \sum_{i=1}^n \frac{1}{2\tau^2} \|G(\mathbf{m}_i) - \tilde{\mathbf{x}}_i^T\|_2^2$
  - 7:   **if** using ERM **then**
  - 8:     Sample  $\tilde{\mathbf{m}}_i \sim q_\alpha(\mathbf{m}|\mathbf{x})$  to get sample pairs  $\{\mathbf{x}_i, z_i, \tilde{\mathbf{m}}_i\}_{i=1}^n$
  - 9:     Compute  $L_{ELBO} = \frac{1}{n} \sum_{i=1}^n [\log p_g(\mathbf{x}_i, z_i|\mathbf{m}_i) + \text{KL}(q_\alpha(\mathbf{m}|\mathbf{x}_i)\|p(\mathbf{m}))]$  via Eq.13.
  - 10:    Compute  $L_G = L_G + L_{ELBO}$
  - 11:   **end if**
  - 12:   Update  $G, q_\alpha$  (if using ERM) by minimizing  $L_G$   
     ▷ EBM training
  - 13:   Apply data augmentation to generated samples:  $\mathbf{x}_{\text{aug}_i} = v(G(\mathbf{m}_i))$ ,  $v \sim \mathcal{V}$  with probability  $p$  and  $G(\mathbf{m}_i)$  otherwise
  - 14:   Sample negative examples  $\{\tilde{\mathbf{x}}_i^T\}_{i=1}^n \sim p_\theta(\mathbf{x})$  initialized from  $\tilde{\mathbf{x}}_i^0 = \mathbf{x}_{\text{aug}_i}$  using Eq.9.
  - 15:   Compute  $L_{\text{EBM}} \leftarrow \frac{1}{n} \sum_{i=1}^n E_\theta(\mathbf{x}_i, z_i) - E_\theta(\tilde{\mathbf{x}}_i^T)$
  - 16:   Update  $E_\theta$  by minimizing  $\mathcal{L}_{\text{EBM}}$
  - 17: **end for**
- 

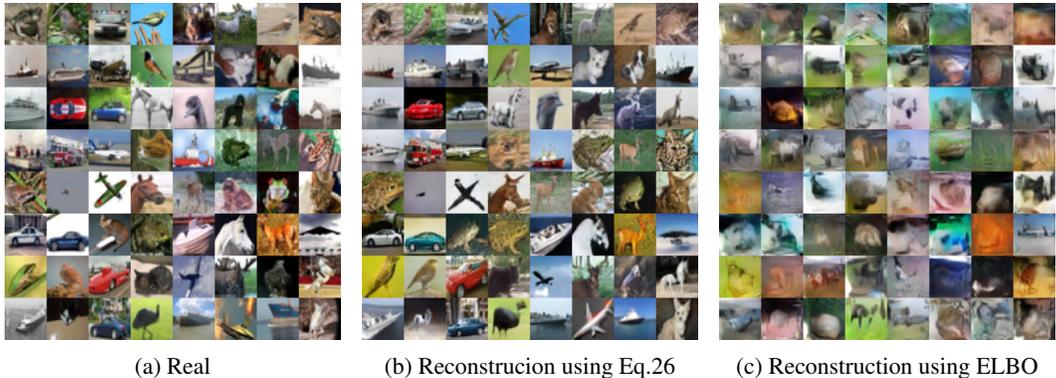
## A.6 INFERENCE MODEL FOR EM SETTING

We train an inference model for the EM setting using the following loss:

$$-\mathbb{E}_{p_{\text{data}}(\mathbf{x}, z)} \mathbb{E}_{q_\alpha(\mathbf{m}|\mathbf{x})} \left[ \log p_g(z|\mathbf{m}) - \frac{q_\alpha(\mathbf{m}|\mathbf{x})}{p(\mathbf{m})} \right], \quad (26)$$

where  $p_g(z|\mathbf{m}) = \rho \text{sim}(z, h(G(\mathbf{m})))$  by definition in Eq.13. The loss essentially minimizes the KL divergence between two conditional distributions:  $\text{KL}(p_{\text{data}}(z|\mathbf{x})\|p_g(z|\mathbf{x}))$ . This approach ensures feature preservation in the latent space rather than enforcing pixel-level reconstruction. Fig.9 compares reconstruction results using Eq.26 versus the traditional ELBO loss in VAEs. The traditional

918 ELBO fails to produce clear, semantically meaningful images with recognizable objects. Our train-  
 919 ing loss achieves high-quality reconstructions that preserve semantic properties of the input, such as  
 920 object class, color, and visual style, without enforcing exact image reproduction. This indicates that  
 921 our latent representation supports flexible instance generation.  
 922



934 (a) Real (b) Reconstruction using Eq.26 (c) Reconstruction using ELBO  
 935  
 936 Figure 9: Reconstruction with different training losses.  
 937

938  
 939 A.7 RECONSTRUCTION OF LIEBM-ERM

940 While our autoencoder-style ERM scheme is designed primarily for initialization, we additionally  
 941 demonstrate its image reconstruction capabilities in Figs.10 and 11. Following the test setting in Han  
 942 et al. (2019), we also compare our approach with other models that also incorporate an inferential  
 943 mechanism, where performance is quantitatively measured by per-pixel mean square error (MSE).  
 944 As shown in Tab.11, our model achieves the best performance on CIFAR-10, outperforming Dual-  
 945 MCMC even with Langevin refinement. On CelebA 64, our model achieves comparable results to  
 946 Dual-MCMC but without requiring additional Langevin dynamics.  
 947

948 Table 11: Reconstruction evaluation using MSE on CIFAR-10 and CelebA 64. Inf+L=10 denotes  
 949 using 10-step Langevin dynamics initialized by the inference model.

Methods	CIFAR-10	CelebA-64
WS (Hinton et al., 1995)	0.058	0.152
VAE (Kingma & Welling, 2014)	0.037	0.039
ALI (Dumoulin et al., 2016)	0.311	0.519
ALICE (Li et al., 2017)	0.034	0.046
Divergence Triangle (Han et al., 2019)	0.028	0.030
Dual-MCMC (Inf) (Cui & Han, 2023)	0.049	0.022
Dual-MCMC (Inf+L=10) (Cui & Han, 2023)	0.024	<b>0.013</b>
<b>LIEBM-ERM (Inf)</b>	<b>0.019</b>	0.014

960  
 961  
 962 A.8 HYPERPARAMETER SETTINGS

963 We specify the hyperparameters used for our training on each dataset in Tab.12. We adopt two  
 964 forms of function  $F$  in Eq.5 for different datasets based on the generation performance. For CIFAR-  
 965 10 and ImageNet 32, we define  $F(f_\phi(x)) = \frac{-\|f_\phi(x)\|_2^2}{2}$ , while for CelebA 64 and CelebA-HQ 256,  
 966 we define  $F$  to be a learnable linear function, which is trained along with  $E_\theta(x, z)$ . The output  
 967 dimension of  $f_\phi(x)$  is 512.  
 968

969  
 970 A.9 ADDITIONAL RESULTS

971 We provide more qualitative visual results for both EM and ERM settings in Figs.12-15.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025



Figure 10: Reconstruction of LIEBM-ERM on CIFAR-10.

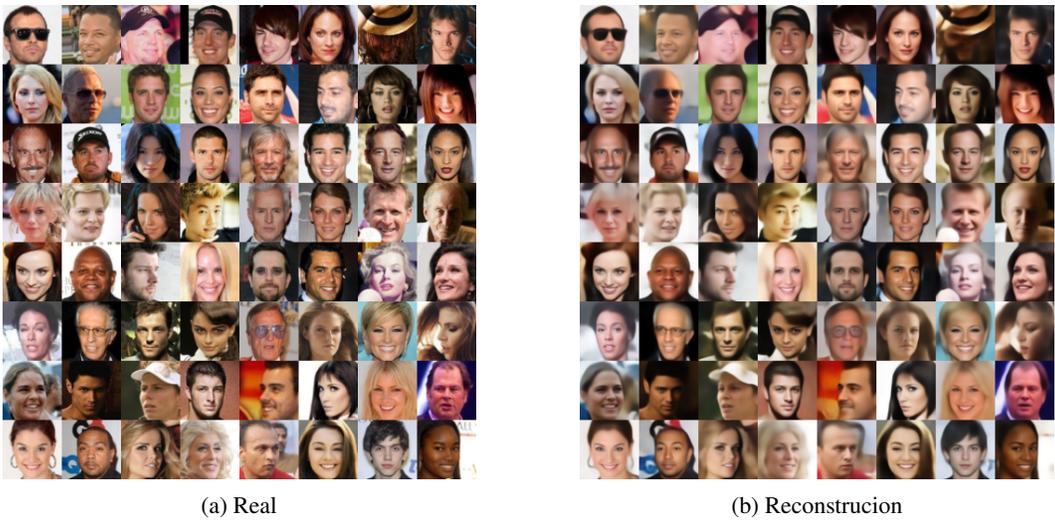
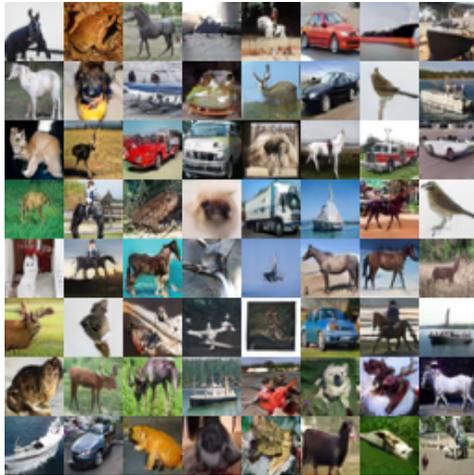


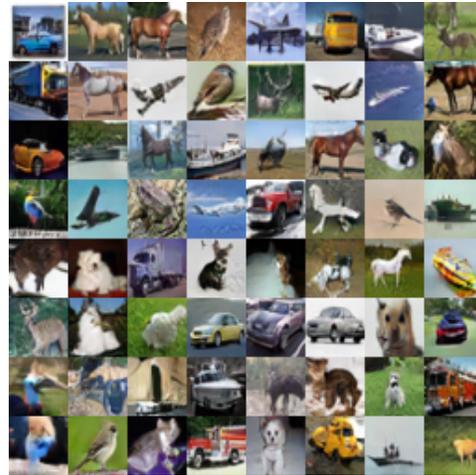
Figure 11: Reconstruction of LIEBM-ERM on CelebA 64.

Table 12: Hyperparameters for each dataset.

	CIFAR-10	ImageNet 32	CelebA 64	CelebA-HQ 256
$E_\theta$ learning rate / Adam $\beta_1, \beta_2$	1e-4 / (0.0, 0.999)	1e-4 / (0.0, 0.999)	1e-4 / (0.0, 0.9)	1e-4 / (0.0, 0.9)
$G$ learning rate / Adam $\beta_1, \beta_2$	2e-4 / (0.0, 0.9)	2e-4 / (0.0, 0.9)	3e-4 / (0.0, 0.9)	3e-4 / (0.0, 0.9)
$q_\alpha$ learning rate / Adam $\beta_1, \beta_2$	2e-4 / (0.0, 0.9)	2e-4 / (0.0, 0.9)	1e-4 / (0.0, 0.9)	1e-4 / (0.0, 0.9)
EMA decay rate	0.9999	0.9999	0.9999	0.9999
$\gamma$ for training	0.01	0.01	0.01	0.01
$\gamma$ for OOD	0.1	0.1	0.1	1
batch size	256	256	256	128
MCMC steps	15	15	15	15
MCMC step size $\delta^2$	25	25	0.1	0.1
$\omega_1 / \omega_2$ in Eq.11	1 / 0.1	1 / 0.1	70 / 1	70 / 1
$\rho$ in Eq.13	1	1	1	50
training epochs	500	100	300	300
data range	[0, 1]	[-1, 1]	[-1, 1]	[-1, 1]
latent dimension	128	128	128	256
$E_\theta, G$ hidden channels	256	512	1024	1024
$q_\alpha$ hidden channels	128	128	128	64
$G$ params	4.3M	16.0M	12.2M	34.3M
$E_\theta$ params	4.9M	17.6M	20.7M	40.7M
$q_\alpha$ params	15.2M	15.2M	15.2M	8.1M



(a) EM



(b) ERM

Figure 12: Samples generated by LIEBM with MCMC refinement on CIFAR-10.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098

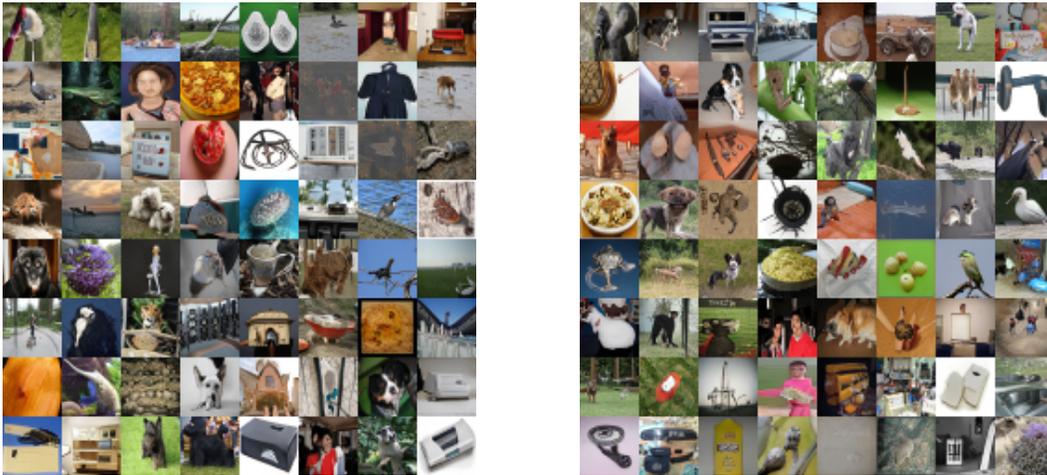


Figure 13: Samples generated by LIEBM with MCMC refinement on ImageNet 32.

1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117



(a) EM

1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129



(b) ERM

Figure 14: Samples generated by LIEBM with MCMC refinement on CelebA 64.

1130  
1131  
1132  
1133

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187



(a) EM



(b) ERM

Figure 15: Samples generated by LIEBM with MCMC refinement on CelebA-HQ 256.