

M-EBM: Towards Understanding the Manifolds of Energy-Based Models

Xiulong Yang^[0000-0003-3417-7106] and Shihao Ji

Georgia State University
xyang22@gsu.edu, sji@gsu.edu

Abstract. Energy-based models (EBMs) exhibit a variety of desirable properties in predictive tasks, such as generality, simplicity and compositionality. However, training EBMs on high-dimensional datasets remains unstable and expensive. In this paper, we present a Manifold EBM (M-EBM) to boost the overall performance of unconditional EBM and Joint Energy-based Model (JEM). Despite its simplicity, M-EBM significantly improves unconditional EBMs in training stability and speed on a host of benchmark datasets, such as CIFAR10, CIFAR100, CelebA-HQ, and ImageNet 32x32. Once class labels are available, label-incorporated M-EBM (M-JEM) further surpasses M-EBM in image generation quality with an over 40% FID improvement, while enjoying improved accuracy. The code can be found in <https://github.com/sndnyang/mebm>.

Keywords: Generative Model · Energy-based Model · Joint Energy-based Model.

1 Introduction

Energy-Based Models (EBMs) are an class of probabilistic models, which are widely applicable in image generation, out of distribution detection, adversarial robustness, and hybrid discriminative-generative modeling [16,4,3,5,7,21,6,20]. However, training EBMs on high-dimensional datasets remains very challenging. Most of the works utilize the Markov Chain Monte Carlo (MCMC) sampling [19] to generate samples from the model distribution represented by an EBM. Specifically, they require K -step Langevin Dynamics sampling [19] to generate samples from the model distribution in every iteration, which can be extremely expensive when using a large number of sampling steps, or highly unstable with a small number of steps. The trade-off between the training time and stability prevents the MCMC sampling based EBMs from scaling to large-scale datasets.



Fig. 1. Generated samples of CelebA-HQ 128x128 from our M-EBM.

Recently, there are a flurry of works on improving EBMs. The most recent studies [3,5] on the MCMC-based approach focus on improving the generation quality and stability. However, they still resort to a long sampling chain and requires expensive training. Another branch of works [6,20] augment the EBM with a regularized generator in a GAN-style training to improve the stability and speed, sacrificing the desired property of learning a single object. Moreover, JEM [7] proposes an elegant framework to reinterpret the modern CNN classifier as an EBM and achieves impressive performances in image classification and generation simultaneously. However, it also suffers from the divergence issue of the MCMC-based sampling, and its generative performance falls behind state-of-the-art EBMs. Tackling the limitations of JEM, JEM++ [21] introduces a variety of training procedures and architecture features to improve JEM in terms of accuracy, speed and stability altogether. Furthermore, JEM++ demonstrates a trade-off between classification accuracy and image quality, but it still cannot improve image generation quality notably.

In this paper, we introduce simple yet effective training techniques to improve unconditional EBM and JEM in terms of image generation quality, training stability and speed altogether. First, the informative initialization introduced in JEM++ dramatically improves the training stability and reduces the required MCMC sampling steps. However, it’s not scalable for high-resolution and large-scale datasets. Hence, we introduce a simplified informative initialization that is suitable for unconditional EBM and JEM for high-resolution images and a large number of classes (e.g., 128x128 CelebA-HQ and 1000-class ImageNet 32x32 datasets). We name our models as Manifold EBM (M-EBM) and Manifold JEM (M-JEM) respectively. Second, we find the L_2 regularization of the energy magnitude does not work with the energy function utilized in JEM. To enable L_2 regularization and improve the training stability, we augment the standard softmax classifier with a new energy head, which is then L_2 regularized. Despite the simplicity, these techniques allow us to reduce the number of MCMC sampling steps of EBM dramatically, while retaining or sometimes improving classification accuracy of prior state-of-the-art EBMs.

Our main contributions are summarised as follows:

1. We simplify the informative initialization in JEM++ for the SGLD chain, which stabilizes and accelerates the training of unconditional EBM and JEM, while being scalable for high-resolution and large-scale datasets.
2. Adding an L_2 -regularized energy head on top of a CNN feature extractor to represent an energy function stabilizes the training of JEM. Then we train M-JEM using two mini-batches: one with data augmentation for classification, and the other one without data augmentation for maximum likelihood estimation of EBMs.
3. We conduct extensive experiments on four benchmark datasets. M-EBM matches or outperforms prior state-of-the-art unconditional EBMs, while significantly improves training stability and reduces the number of sampling steps. Moreover, M-JEM improves JEM’s training stability and speed, image

generation quality, and classification accuracy altogether, while outperforming M-EBM in image generation quality.

2 Background

Energy-based Models (EBMs) [13] utilizes the idea that any probability density $p_{\theta}(\mathbf{x})$ can be expressed as

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{Z(\theta)}, \quad (1)$$

where $E_{\theta}(\mathbf{x})$ is named the energy function that maps each input $\mathbf{x} \in \mathcal{X}$ to a scalar, and $Z(\theta) = \int_{\mathbf{x}} \exp(-E_{\theta}(\mathbf{x})) d\mathbf{x}$ is the normalizing constant w.r.t \mathbf{x} (also known as the partition function). Ideally, an energy function should assign low energy values to samples drawn from data distribution and high values otherwise.

The key challenge of EBM training is estimating the intractable partition function $Z(\theta)$, and the maximum likelihood estimation of parameters θ is not straightforward. A number of sampling-based approaches have been proposed to approximate the partition function effectively. Specifically, the derivative of the log-likelihood of $\mathbf{x} \in \mathcal{X}$ w.r.t. θ can be expressed as

$$\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[\frac{\partial E_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \mathbb{E}_{p_d(\mathbf{x})} \left[\frac{\partial E_{\theta}(\mathbf{x})}{\partial \theta} \right], \quad (2)$$

where the first expectation is over the model density $p_{\theta}(\mathbf{x}')$, which is challenging due to the intractable $Z(\theta)$.

To estimate it efficiently, MCMC and Gibbs sampling [10] have been proposed. Moreover, to speed up the sampling, recently Stochastic Gradient Langevin Dynamics (SGLD) [19] is employed to train EBMs [16,4,7]. Specifically, to sample from $p_{\theta}(\mathbf{x})$, the SGLD follows

$$\mathbf{x}^0 \sim p_0(\mathbf{x}), \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\alpha}{2} \frac{\partial E_{\theta}(\mathbf{x}^t)}{\partial \mathbf{x}^t} + \alpha \epsilon^t, \quad \epsilon^t \sim \mathcal{N}(0, 1), \quad (3)$$

where $p_0(\mathbf{x})$ is typically a uniform distribution over $[-1, 1]$, whose samples are refined via a noisy gradient decent with step-size α over a sampling chain.

Prior works [15,16,4,7] have investigated the effect of hyper-parameters in SGLD sampling and showed that the SGLD-based approaches suffer from poor stability and prolonged computation of sampling at every iteration. Nijkamp et al. [15] find that it's desirable to generate samples from the SGLD chain after it converges. The convergence requires the step-size α to decay with a polynomial schedule and infinite sampling steps, which is impractical. Therefore, Short-Run and Long-Run MCMC samplings are utilized for EBM training. Moreover, most works [4,7,3] use a constant step-size α during sampling and approximate the samples with a sampler that runs only for a finite number of steps, which is still computationally very expensive. Another recent work [5] combines the

SGLD-based approach with diffusion models [11] under a framework of conditional EBMs. They achieve state-of-the-art image generation quality and obtain a faithful energy potential.

Joint Energy-based Models (JEM) [7] demonstrates that standard softmax-based classifiers can be trained as EBMs. Given an input $\mathbf{x} \in R^D$, a classifier of parameters θ maps the input to a vector of C real-valued numbers (known as logits): $f_\theta(\mathbf{x})[y], \forall y \in [1, \dots, C]$, where C is the number of classes. Then the softmax function is employed to convert the logits into a categorical distribution: $p_\theta(y|\mathbf{x}) = e^{f_\theta(\mathbf{x})[y]} / \sum_{y'} e^{f_\theta(\mathbf{x})[y']}$. The authors reuse the logits to define an energy function for the joint density: $p_\theta(\mathbf{x}, y) = e^{f_\theta(\mathbf{x})[y]} / Z(\theta)$. Then a marginal density of \mathbf{x} can be achieved by marginalizing out y as: $p_\theta(\mathbf{x}) = \sum_y p_\theta(\mathbf{x}, y) = \sum_y e^{f_\theta(\mathbf{x})[y]} / Z(\theta)$. As a result, the corresponding energy function of \mathbf{x} is defined as

$$E_\theta(\mathbf{x}) = -\log \sum_y e^{f_\theta(\mathbf{x})[y]} = -\text{LSE}(f_\theta(\mathbf{x})), \quad (4)$$

where $\text{LSE}(\cdot)$ denotes the Log-Sum-Exp function. The advantage of this LSE energy function is that an additional degree of freedom in the scale of the logit vector now can model the data distribution.

To optimize the model parameter θ , JEM maximizes the logarithm of joint density function $p_\theta(\mathbf{x}, y)$:

$$\log p_\theta(\mathbf{x}, y) = \log p_\theta(y|\mathbf{x}) + \log p_\theta(\mathbf{x}), \quad (5)$$

where the first term is the cross-entropy objective for classification, and the second term is the maximum likelihood learning of EBM as shown in Eq. 2. We can also interpret the second term as an unsupervised regularization on the model parameters θ .

3 Manifold EBM

3.1 Informative Initialization and M-EBM

As shown in Eq. 3, the SGLD sampling starts from an initial distribution $p_0(\mathbf{x})$. To train the EBM as a generative model, Short-Run MCMC sampling [16] utilizes an MCMC sampler that starts from a random noise distribution such as a uniform distribution. A concurrent work IGEBM [4] proposes an initialization approach with a sample replay buffer in which they store past generated samples and draw samples from either replay buffer or uniform random noise to initialize the Langevin dynamics procedure. This is also the sampling approach adopted by [7,25]. Furthermore, JEM++ [21] introduces an informative initialization with the replay buffer by using a Gaussian mixture distribution estimated from the training images, which significantly reduces the number of sampling steps required by SGLD while improving its training stability.

However, the per-class covariance matrices of the Gaussian mixture distribution utilized by JEM++ can be huge for high-resolution image datasets with a

large number of classes. Hence, we estimate a single Gaussian distribution from the whole training dataset. That is, we estimate the initial sampling distribution as

$$p_0(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (6)$$

with $\boldsymbol{\mu} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}]$, $\boldsymbol{\Sigma} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]$,

where \mathcal{D} denotes the whole training set. The visualization of the estimated centers and samples from $p_0(\mathbf{x})$ of different datasets are provided in the appendix. Since only one Gaussian distribution is estimated from the whole training set, we can apply it for unconditional datasets such as CelebA, and reduce the memory and space required for the large covariance matrices¹. Although $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ can be well estimated with sufficient samples, they still lead to a biased initialization with higher variance, compared to the Gaussian mixture initialization utilized in JEM++. But our empirical study shows that our simplified initialization won't deteriorate the performance and is comparable to bias-reduced Gaussian mixture initialization.

Since the manifold of \mathbf{x}_0 from our informative initialization is much closer to the real data manifold than that of uniform initialization, this informative initialization reduces the required sampling steps (and thus accelerates training), and also improves training stability as we will demonstrate in the experiments. We therefore call the EBM with this simplified informative initialization as M-EBM throughout this work.

3.2 Injected Noise in M-EBM

Existing work [16] studied the effect of injected noise on training stability via smoothing p_{data} with additive Gaussian noises $\mathbf{x} \leftarrow \mathbf{x} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 I)$. Their results demonstrated that the fidelity of the examples in terms of IS and FID improves, when lowering σ^2 . And they depict the tradeoff between the sampling steps K and the level of injected noise, indicating the training time and the stability. After it, several following methods [5,3,22] successfully remove the injected noise and achieve better image quality. However, they require a very large $K \geq 30$ to stabilize the training. Thanks to the informative initialization, it not only allows us to significantly reduce K , but also removes the injected noise to improve the image quality while keeping high stability. As shown in Fig. 4(b), the manifolds of real data and \mathbf{x}_0 sampled from informative initialization are very close, even mixing together when M-EBM is trained without energy regularization. Hence, we suppose the gradients $\nabla_{\mathbf{x}} E(\mathbf{x})$ are defined (almost) everywhere in such manifolds and thus can reduce the perturbation with noise which is originally explained in NCSN [18].

¹ One covariance matrix of CIFAR10 has $(3 \times 32 \times 32)^2 \approx 9.4M$ parameters and uses 37.6MB memory. A dataset with C classes will take $37.6 \times C$ MB.

4 Manifold JEM

4.1 Injected Noise in M-JEM

As discussed in previous section, the injected noise smoothing p_{data} would hurt the generative performance of EBMs. For JEM and JEM++, we suppose it would also decrease the classification accuracy. Hence, it’s critical to remove the injected noise and gain benefits in terms of classification accuracy and generation quality. We use N to denote the noise-adding operation. Then the actual objective of JEM is

$$\log p_{\theta}(N(\mathbf{x}), y) = \log p_{\theta}(y|N(\mathbf{x})) + \log p_{\theta}(N(\mathbf{x})). \quad (7)$$

Interestingly, we find that if we only remove the injected noise, the training is not stable. However, if we further disable the data augmentation when learning maximum likelihood $\log p_{\theta}(\mathbf{x})$, it becomes even more stable than JEM++ and enjoys improved accuracy and better sampling quality. Following the observation, we train our M-JEM using two mini-batches: one with data augmentation for classification, and the other one without data augmentation for maximum likelihood estimation of EBMs.

4.2 Energy Function Regularization in M-JEM

IGEBM [4] finds that constraining the Lipschitz constant of the energy network can ease the instability issue in Langevin dynamics. Hence, they weakly L_2 regularize energy magnitudes for both positive and negative samples to the contrastive divergence as:

$$\mathcal{L} = \frac{1}{B} \sum_{i=1}^B \left(E_i^+ - E_i^- + \alpha(E_i^{+2} + E_i^{-2}) \right), \quad (8)$$

where $E^+ = E_{\theta}(\mathbf{x}^+)$ with \mathbf{x}^+ sampled from the data distribution p_d , and $E^- = E_{\theta}(\mathbf{x}^-)$ with \mathbf{x}^- sampled from the model distribution $p_{\theta}(\mathbf{x})$. The effect of L_2 regularization on EBMs can be viewed as Fig 4(b). However, since L_2 regularization would force the vector of logits $f_{\theta}(\mathbf{x})$ to be uniform, while maximizing $p_{\theta}(y|\mathbf{x})$ boosts $f_{\theta}(\mathbf{x})[y]$. Hence, the L_2 regularization is incompatible with Eq. 4 and cannot be directly applied to vanilla JEM.

To incorporate L_2 regularization to JEM, we propose to augment the standard CNN softmax classifier with an extra fully connected layer, called Energy Head, as shown in Fig. 2(a). Then the L_2 regularization is applied on the energy head (instead of the LSE classification head) to improve the training stability.

We provide the pseudo-code for M-EBM/JEM as in Algorithm 1, which follows the framework of IGEBM [4] and JEM [7].

5 Experiments

In this section, we first evaluate the generative performance of M-EBM on multiple datasets, including CIFAR10, CIFAR100, CelebA-HQ 128x128 and ImageNet

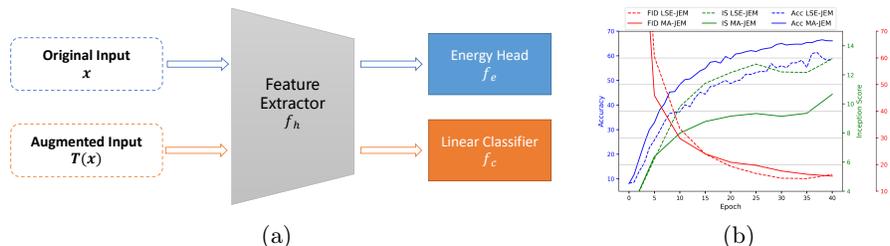


Fig. 2. a) The architecture of M-JEM. An energy head f_e is augmented for energy magnitude regularization and two mini-batches are used for the training of classifier and the maximum likelihood estimate of EBM, respectively. b) Comparison between M-JEM and LSE-JEM on CIFAR100.

Algorithm 1 M-EBM/JEM Training: Given network f_θ , SGLD step-size α , SGLD noise σ , replay buffer B , SGLD steps K , reinitialization frequency ρ

- 1: **while** not converged **do**
 - 2: Sample \mathbf{x}^+ and y from dataset
 - 3: Sample $\hat{\mathbf{x}}_0 \sim B$ with probability $1 - \rho$, else $\hat{\mathbf{x}}_0 \sim p_\theta(\mathbf{x})$ as Eq. 6
 - 4: **for** $t \in [1, 2, \dots, K]$ **do**
 - 5: $\hat{\mathbf{x}}_t = \hat{\mathbf{x}}_{t-1} - \alpha \cdot \frac{\partial E(\hat{\mathbf{x}}_{t-1})}{\partial \hat{\mathbf{x}}_{t-1}} + \sigma \cdot \mathcal{N}(0, I)$
 - 6: **end for**
 - 7: $\mathbf{x}^- = \text{StopGrad}(\hat{\mathbf{x}}_K)$
 - 8: $L_{\text{gen}}(\theta) = E(\mathbf{x}^+) - E(\mathbf{x}^-) + \alpha (E(\mathbf{x}^+)^2 + E(\mathbf{x}^-)^2)$ as Eq. 8.
 - 9: $L(\theta) = L_{\text{gen}}(\theta)$ for M-EBM
 - 10: $L(\theta) = L_{\text{clf}}(\theta) + L_{\text{gen}}(\theta)$ with $L_{\text{clf}}(\theta) = \text{xent}(f_\theta(\mathbf{x}), y)$ for M-JEM
 - 11: Calculate gradient $\frac{\partial L(\theta)}{\partial \theta}$ to update θ
 - 12: Add \mathbf{x}^- to B
 - 13: **end while**
-

32x32. Then, we investigate the efficacy of the M-JEM on CIFAR10 and CIFAR100. Finally, the study and visualization of the differences between trained EBM and JEM are provided to analyze their generative capability.

Our code is largely built on top of JEM [7]². For a fair comparison with JEM, we update each model with 390 iterations in 1 epoch. Empirically, we find a batch size of 128 for $p_\theta(y|\mathbf{x})$ achieves the best classification accuracy on CIFAR10, while we use 64, the same batch size as in JEM, for $p_\theta(\mathbf{x})$. We train our models on ImageNet 32x32 for 50 epochs and other datasets for 150 epochs at most. All our experiments are performed with PyTorch on Nvidia GPUs. For CIFAR10 and CIFAR100, we train the backbone Wide-ResNet 28-10 [23] on a single GPU. Due to limited computational resources, we use Wide-ResNet 28-2 for ImageNet 32x32 on a single GPU, and Wide-ResNet 28-5 for CelebA-HQ 128x128 on 2 GPUs.

² <https://github.com/wgrathwohl/JEM>

Table 1. Inception and FID scores of M-EBM on CIFAR10.

Model	IS \uparrow	FID \downarrow
M-EBM(K=1)*	6.02	35.7
M-EBM(K=2)	6.72	27.1
M-EBM(K=5)	7.14	22.7
M-EBM(K=10)	7.08	20.4
M-EBM(K=20)	7.20	21.1
Explicit EBM(Unconditional)		
ShortRun(K=100) [16]	6.72	32.1
IGEBM(K=60) [4]	6.78	38.2
f-EBM(K=60) [22]	8.61	30.8
CF-EBM(K=50) [24]	-	16.7
KL-EBM(K=40) [3]	7.85	25.1
DiffuRecov(K=30) [5]	8.31	9.58
Regularized Generator		
GEBM [1]	-	23.02
VAEBM(K=6) [20]	8.43	12.19
Other		
SNGAN [14]	8.59	21.7
NCSN [18]	8.91	25.3
StyleGAN2-ADA [12]	9.74	2.92
DDPM [11]	9.46	3.17

* M-EBM diverges with $K = 1$, and we report the best FID before diverging.

Table 2. FID results of M-EBM on CIFAR100, CelebA-HQ 128, and ImageNet 32x32.

Model	FID \downarrow
CIFAR100 Unconditional	
M-EBM(K=1)*	45.5
M-EBM(K=2)	26.2
M-EBM(K=5)	27.2
M-EBM(K=10)	26.9
SNGAN [14]	22.4
CelebA-HQ 128 Unconditional	
M-EBM(K=5)*	57.76
M-EBM(K=10)	39.87
KL-EBM(K=40) [3]	28.78
SNGAN [14]	24.36
ImageNet 32x32 Unconditional	
M-EBM(K=2)	54.52
M-EBM(K=5)	52.71
IGEBM(K=60) [4]	62.23
KL-EBM(K=40) [3]	32.48

* Our models diverge during training with given K , and we report the best FID before diverging.

5.1 M-EBM

We first evaluate the performance of M-EBM on CIFAR10, CIFAR100, CelebA-HQ 128 and ImageNet 32x32. We utilize the Inception Score (IS) [17] and Fréchet Inception Distance (FID) [9] to evaluate the quality of generated images.

The results are reported in Table 1 and 2, respectively. It can be observed that our method consistently surpasses existing methods in terms of sampling steps by a significant margin. On CIFAR10, M-EBM outperforms many EBM approaches and SNGAN in terms of FID, while the performance is slightly worse than SNGAN on CIFAR100. Some EBM approaches show better performance, such as VAEBM, CF-EBM and DiffuRecov. However, they require an extra pre-trained generator, or special architecture, or much larger sampling steps, while M-EBM can train on a classical architecture as the backbone with least K . On ImageNet 32x32, we note that M-EBM with $K = 2$ is incredibly stable and achieves FID 54.52 within 30 epochs and outperforms IGEBM. In addition, increasing sampling steps K further doesn't have an obvious improvement. Finally, on CelebA-HQ, M-EBM is worse than baseline methods as we find it

is less stable and requires more sampling steps due to the high resolution of CelebA-HQ. Nevertheless, our method builds a new solid baseline on different large-scale benchmarks for further investigations of EBM training in these more challenging tasks. Samples generated by M-EBMs for CIFAR10, CIFAR100, and CelebA-HQ are shown in Fig. 3 and Fig. 1, respectively. The generated samples of ImageNet 32x32 can be found in the appendix.

5.2 M-JEM

We train M-JEM on two benchmark datasets: CIFAR10 and CIFAR100, and compare its performance to the state-of-the-art hybrid models and some representative generative models. Table 3 and 4 report results on CIFAR10 and CIFAR100, respectively. As we can see, M-JEM improves JEM’s image generation quality, stability, speed, and accuracy by a notable margin. It also boosts the IS and FID scores over M-EBM. Compared with JEM++, FID of M-JEM drops dramatically since we exclude the noise, and the notable gain of accuracy when $K = 5$ indicates M-JEM($K = 5$) is much more stable than JEM++($K = 5$). On CIFAR100, IS and FID scores are not commonly reported by state-of-the-art hybrid models, such as JEM [7], VERA [6], and JEM++ [21]. Hence, our work builds a baseline for hybrid modeling on CIFAR100 with decent classification accuracy and image generation quality for future investigations. Images generated by M-JEM for CIFAR10 and CIFAR100 are can be found in Fig. 3.

5.3 Analysis

Is Energy Head better than LSE? To evaluate the effect of the energy head, we conduct an experiment comparing M-JEM (with energy head) and LSE-JEM (without energy head) on CIFAR100. Fig. 2(b) shows that M-JEM achieves much higher classification accuracy, comparable FID but a lower Inception Score than LSE-JEM. However, we empirically find LSE-JEM is less stable than M-JEM after 40 epochs which leads us to analyze the manifolds learned by different models.

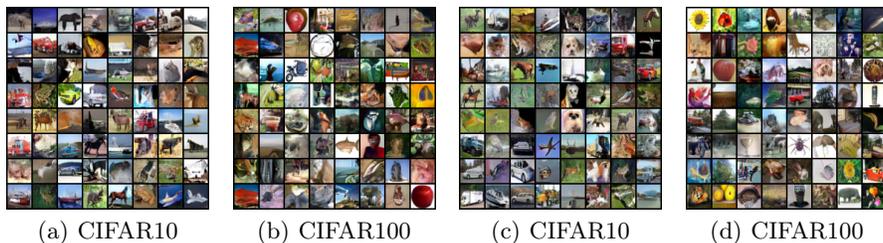


Fig. 3. M-EBM and M-JEM generated samples of CIFAR10 and CIFAR100.

Manifold Analysis To facilitate better understanding of different approaches, we utilize the t-SNE visualization for manifold analysis shown in Fig. 4(a) and 4(b). To have a fair comparison, we pick fixed samples from CIFAR10 as \mathbf{x}^+ , initialize samples from $p_0(\mathbf{x})$ as \mathbf{x}^0 , and randomly select samples from the replay buffer of each pre-trained models as \mathbf{x}^- . Given the inputs from \mathbf{x}^+ , \mathbf{x}^- and \mathbf{x}^0 , we collect the outputs of the penultimate layer as features and apply the t-SNE technique to generate the visualization. For Fig 4(a), three CIFAR10-trained M-EBM, M-JEM, and LSE-JEM with $K = 10$ are involved. We further conduct the comparison between M-EBMs($K = 5$) with and without energy L_2 regularization in Fig 4(b).

Table 3. Hybrid Modeling Results on CIFAR10.

Model	Acc % \uparrow	IS \uparrow	FID \downarrow
M-JEM(K=1)*	78.4	7.91	29.8
M-JEM(K=2)*	86.5	8.64	19.3
M-JEM(K=5)	93.1	8.71	12.1
M-JEM(K=10)	93.8	8.52	11.5
M-JEM(K=20)	94.2	8.72	12.2

Model	Acc % \uparrow	IS \uparrow	FID \downarrow
Softmax	78.9	-	-
SNGAN(Cond)	-	9.30	15.6
BigGAN(Cond)	-	11.0	11.73

Model	Acc % \uparrow	IS \uparrow	FID \downarrow
Residual Flow [2]	70.3	3.60	46.4
IGEBM(K=60) [4]	49.1	8.30	37.9
JEM(K=20) ⁺ [7]	92.9	8.76	38.4
JEM++(M=5) ⁺ [21]	91.1	7.81	37.9
JEM++(M=10)	93.5	8.29	37.1
JEM++(M=20)	94.1	8.11	38.0
JEAT [25]	85.2	8.80	38.2
M-JEM(K=20)*	70.4	10.32	51.7
JEM(K=30)*	72.8	10.84	34.2
JEM++(K=5)*	72.0	8.19	37.7
JEM++(K=10)*	74.5	10.23	32.9
VERA($\alpha=100$)*	69.3	8.14	28.2
VERA($\alpha=1$)*	48.7	7.97	26.6
M-JEM(K=1) ⁺	46.5	8.71	26.2
M-JEM(K=2) ⁺	63.5	11.22	15.1
M-JEM(K=5)	73.5	11.95	13.5
M-JEM(K=10)	75.1	11.72	12.7

VERA($\alpha=100$)	93.2	8.11	30.5
VERA($\alpha=1$) [6]	76.1	8.00	27.5
softmax	95.8	-	-

* We report the best performance before the diverging of training.

⁺ They suffer from high instability and regularly diverge.

As we can observe in Fig. 4(a), M-JEM with label information forms more compact manifolds of \mathbf{x}^+ and \mathbf{x}^- than M-EBM. In other words, M-JEM-generated samples \mathbf{x}^- match the distribution of real data and have lower variance and less **manifold intrusion**[8] than M-EBM. It gives us an explanation of why label information can improve generation quality. Moreover, the latent feature space of M-JEM is better formulated than LSE-JEM, and there's less overlap between \mathbf{x}^0 and \mathbf{x}^+ which is desired since \mathbf{x}^0 and \mathbf{x}^+ should be assigned with different energies. Intuitively, the number K of SGLD sampling required for stable training is correlated to the distance between manifolds of \mathbf{x}^0 and \mathbf{x}^+ . However, in Figures 4(a) and 4(b), we can also observe that \mathbf{x}^0 and \mathbf{x}^+ are roughly mixed

Table 4. Hybrid Modeling Results on CIFAR100.

* No official IS and FID scores are reported.
[‡] We report the best FID before diverging.

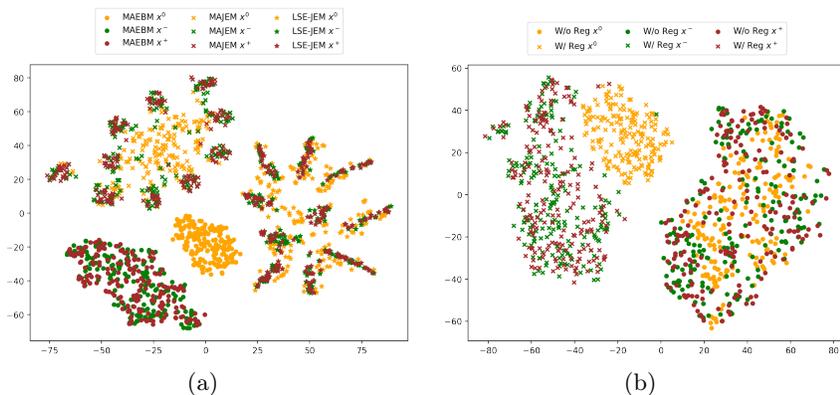


Fig. 4. t-SNE visualization of the latent feature spaces learned by different models trained on CIFAR10. We use different colors to represent (x^0 initial samples, x^- samples in replay buffer, x^+ real data), and different shapes(\circ, \times, \star) to indicate different EBMs.

together from M-EBM without regularization and LSE-JEM. Hence, it’s interesting to reconsider the distance and the training instability when x^0 and x^+ are somewhat mixing together. We leave the exploration of this phenomenon as an existing direction for future.

6 Conclusion

In this paper, we propose simple yet effective training techniques to improve the image generation quality, training speed, and stability of unconditional EBM and JEM altogether. The experimental results demonstrate that our models achieve comparable performance on unconditional EBMs and JEMs, and enable us to scale the MCMC-based EBM learning to high-resolution large-scale image datasets, such as CelebA-HQ 128x128 and ImageNet 32x32 with the least MCMC sampling steps, making EBM training more practical for a research lab in academia to afford and explore.

References

1. Arbel, M., Zhou, Liang and Gretton, A.: Generalized energy based models. In: International Conference on Learning Representations (ICLR) (2021)
2. Chen, R.T., Behrmann, J., Duvenaud, D., Jacobsen, J.H.: Residual flows for invertible generative modeling. In: Neural Information Processing Systems (2019)
3. Du, Y., Li, S., Tenenbaum, J., Mordatch, I.: Improved Contrastive Divergence Training of Energy Based Models. In: International Conference on Machine Learning (ICML) (2021)
4. Du, Y., Mordatch, I.: Implicit generation and generalization in energy-based models. In: Neural Information Processing Systems (NeurIPS) (2019)
5. Gao, R., Song, Y., Poole, B., Wu, Y.N., Kingma, D.P.: Learning Energy-Based Models by Diffusion Recovery Likelihood. In: International Conference on Learning Representations (ICLR) (2021)

6. Grathwohl, W., Kelly, J., Hashemi, M., Norouzi, M., Swersky, K., Duvenaud, D.: No mcmc for me: Amortized sampling for fast and stable training of energy-based models. In: International Conference on Learning Representations (ICLR) (2021)
7. Grathwohl, W., Wang, K.C., Jacobsen, J.H., Duvenaud, D., Norouzi, M., Swersky, K.: Your classifier is secretly an energy based model and you should treat it like one. In: International Conference on Learning Representations (ICLR) (2020)
8. Guo, H., Mao, Y., Zhang, R.: Mixup as locally linear out-of-manifold regularization. In: AAAI Conference on Artificial Intelligence (AAAI) (2019)
9. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: Neural Information Processing Systems (NeurIPS) (2017)
10. Hinton, G.E.: Training products of experts by minimizing contrastive divergence. *Neural computation* (2002)
11. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: Neural Information Processing Systems (NeurIPS) (2020)
12. Karras, T., Aittala, M., Hellsten, J., Laine, S., Lehtinen, J., Aila, T.: Training generative adversarial networks with limited data. In: Proc. Neural Information Processing Systems (NeurIPS) (2020)
13. LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M., Huang, F.: A tutorial on energy-based learning. *Predicting structured data* (2006)
14. Miyato, T., Kataoka, T., Koyama, M., Yoshida, Y.: Spectral normalization for generative adversarial networks. In: International Conference on Learning Representations (ICLR) (2018)
15. Nijkamp, E., Hill, M., Han, T., Zhu, S.C., Wu, Y.N.: On the anatomy of mcmc-based maximum likelihood learning of energy-based models. In: AAAI Conference on Artificial Intelligence (AAAI) (2019)
16. Nijkamp, E., Zhu, S.C., Wu, Y.N.: Learning non-convergent short-run mcmc toward energy-based model. In: Neural Information Processing Systems (NeurIPS) (2019)
17. Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training gans. In: Neural Information Processing Systems (NeurIPS) (2016)
18. Song, Y., Ermon, S.: Generative modeling by estimating gradients of the data distribution. In: Neural Information Processing Systems (NeurIPS) (2019)
19. Welling, M., Teh, Y.W.: Bayesian learning via stochastic gradient langevin dynamics. In: International Conference on Machine Learning (ICML) (2011)
20. Xiao, Z., Kreis, K., Kautz, J., Vahdat, A.: VAEBM: A Symbiosis between Variational Autoencoders and Energy-based Models. In: International Conference on Learning Representations (ICLR) (2021)
21. Yang, X., Ji, S.: JEM++: Improved Techniques for Training JEM. In: International Conference on Computer Vision (ICCV) (2021)
22. Yu, L., Song, Y., Song, J., Ermon, S.: Training deep energy-based models with f-divergence minimization. In: International Conference on Machine Learning (ICML) (2020)
23. Zagoruyko, S., Komodakis, N.: Wide residual networks. In: BMVC (2016)
24. Zhao, Y., Xie, J., Li, P.: Learning energy-based generative models via coarse-to-fine expanding and sampling. In: International Conference on Learning Representations (ICLR) (2021)
25. Zhu, Y., Ma, J., Sun, J., Chen, Z., Jiang, R., Li, Z.: Towards Understanding the Generative Capability of Adversarially Robust Classifiers. In: IEEE International Conference on Computer Vision (ICCV) (2021)

A Experimental Details

Our code is largely built on top of JEM [7]³. For a fair comparison with JEM, we update each model with 390 iterations in 1 epoch. Empirically, we find a batch size of 128 for $p_{\theta}(y|\mathbf{x})$ achieves the best classification accuracy on CIFAR10, while we use 64, the same batch size as in JEM, for the maximum likelihood estimate of $p_{\theta}(\mathbf{x})$. All our experiments are performed with PyTorch on Nvidia GPUs. For CIFAR10 and CIFAR100, we train the backbone Wide-ResNet 28-10 [23] on a single GPU. Due to limited computational resources, we use Wide-ResNet 28-2 for ImageNet 32x32 on a single GPU, and Wide-ResNet 28-5 for CelebA-HQ 128x128 on 2 GPUs.

Table 5 lists the hyper-parameters of our M-EBM/JEM algorithms. We train all our models for 200 epochs with the SGD optimizer, a buffer size of 10,000, a reinitialization frequency of 5%. For CIFAR10, we use a larger learning rate of 0.1, while for CIFAR100, CelebA-HQ and Imagenet 32x32 we use an initial learning rate of 0.02.

Table 5. Hyper-parameters of M-EBM/JEM

Variable	Value
Epochs N	200
Buffer size $ \mathbb{B} $	10,000
Reinitialization freq. ρ	5%
SGLD step-size α	1
SGLD noise σ	0.001

B Training Speed

We report the empirical training speeds of our M-JEM and baseline methods on a single Titan GPU in Table 6. As discussed previously, two mini-batches are utilized in M-JEM: one for training of EBMs and the other one for training of classifiers. In our experiments, we set the mini-batch size to 64 for EBM training, but use 128 for the classification batch because it achieves the best accuracy. It can be observed that M-EBM/JEM with the least number of MCMC sampling steps (e.g., $K = 2$ and $K = 5$) can train without divergence⁴, and leads to 4.87x and 2.22x speedups over the original JEM, respectively. Here, the acceleration comes from the reduced MCMC sampling steps without stability dropping. M-JEM do not affect the sampling speed since it uses the same backbone and the

³ <https://github.com/wgrathwohl/JEM>

⁴ We can't guarantee the stability. M-EBM($K=2$)/M-JEM($K=5$) sometimes can train without failure, while other EBMs always diverge before 100 epochs

SGLD sampling algorithm as JEM and its variants. Notably, Diffusion-based EBMs [5] and VAEBM [20], requiring 8 GPUs/TPUs for days of training, are too expensive to be practical for a research lab in academia to afford.

Table 6. Run-time comparison on CIFAR10. We set 390 iterations as one epoch and the training epochs are 200 epochs. The batch size is 64 for all models except DiffuRecov and VAEBM.

Model	Minutes per Epochs	Runtime (Hours)	Actual Speedup
Classifier	0.77	2.6	
JEM(K=20)*	15.1	50.3	1×
JEM++(K=5)	6.3	20.9	2.39×
JEM++(K=10)	10.2	33.9	1.48×
VERA	9.6	32.2	2.8× [†]
M-EBM			
K=1	2.4	7.9	6.29×
K=2	3.1	10.3	4.87×
K=5	5.4	18.0	2.79×
K=10	9.0	30.0	1.67×
M-JEM			
K=1	3.8	12.7	3.97×
K=2	4.6	15.3	3.28×
K=5	6.8	22.7	2.22×
K=10	10.5	35.0	1.43×
Other			
IGEEM	1 GPU for 2 days		
KL-EBM	1 GPU for 1 days		
DDPM	800k iter, 8 TPUs, 10.6 hours		
DiffuRecov	240k iterations, 8 TPUs, 40+ hours		
VAEBM ⁺	400 epochs, 8 GPUs, 55 hours		

* JEM(K=20) is much less stable than M-EBM(K=2) and M-JEM(K=5).

[†] VERA reports a 2.8× speedup while we run the official code and report a fair comparison results.

⁺ The runtime is for pretraining NVAE only. For VAEBM, they report the training takes around 25,000 iterations (or 16 epochs) on CIFAR-10 using one 32-GB V100 GPU. Then they cannot generate realistic samples anymore.

C Qualitative Analysis of Samples

Generation quality is difficult to qualify. Following the setting of JEM [7], we conduct a qualitative analysis of samples on CIFAR10.

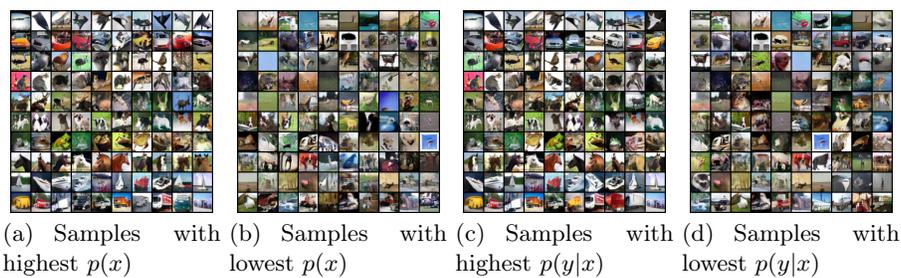


Fig. 5. Each row corresponds to 1 class.

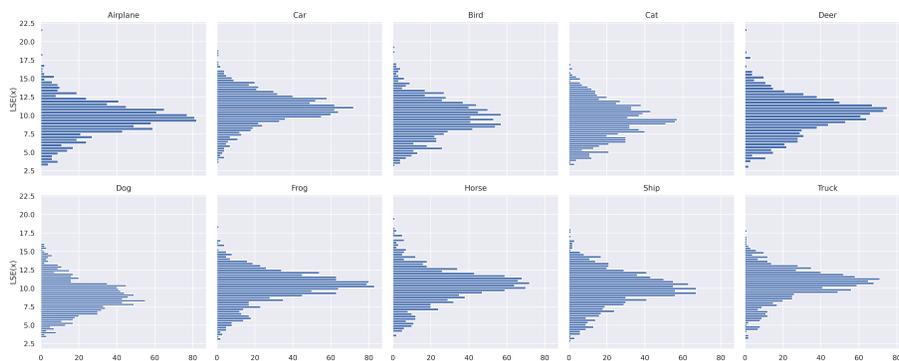


Fig. 6. Histograms (oriented horizontally for easier visual alignment) of $\log p(x)$ arranged by class for CIFAR10.

D Informative Initialization

In this work, we estimate the initial sampling distribution as

$$p_0(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (9)$$

$$\text{with } \boldsymbol{\mu} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x}], \quad \boldsymbol{\Sigma} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \right],$$

where \mathcal{D} denotes the whole training set. Fig. 7 illustrate the center ($\boldsymbol{\mu}$) estimated from each benchmark dataset.

E Additional Generated Samples

Fig. 8 illustrates example images generated by M-EBM on Imagenet 32x32. Additional M-JEM generated class-conditional (best and worst) samples of CIFAR10 are provided in Figures 9- 17.

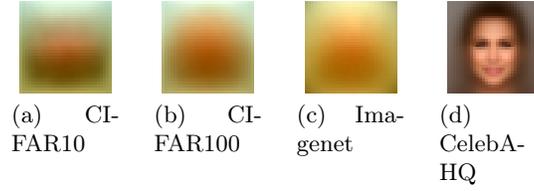
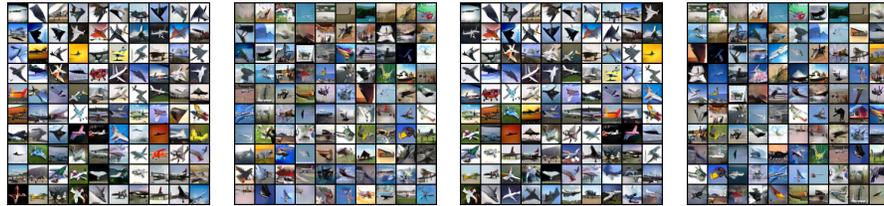


Fig. 7. The centers (μ 's) of CIFAR10, CIFAR100, Imagenet 32 and CelebA-HQ 128.



Fig. 8. M-EBM generated samples of Imagenet 32x32.



(a) Samples with highest $p(\mathbf{x})$ (b) Samples with lowest $p(\mathbf{x})$ (c) Samples with highest $p(y|\mathbf{x})$ (d) Samples with lowest $p(y|\mathbf{x})$

Fig. 9. M-JEM generated class-conditional samples of **Plane**



(a) Samples with highest $p(\mathbf{x})$ (b) Samples with lowest $p(\mathbf{x})$ (c) Samples with highest $p(y|\mathbf{x})$ (d) Samples with lowest $p(y|\mathbf{x})$

Fig. 10. M-JEM generated class-conditional samples of **Car**

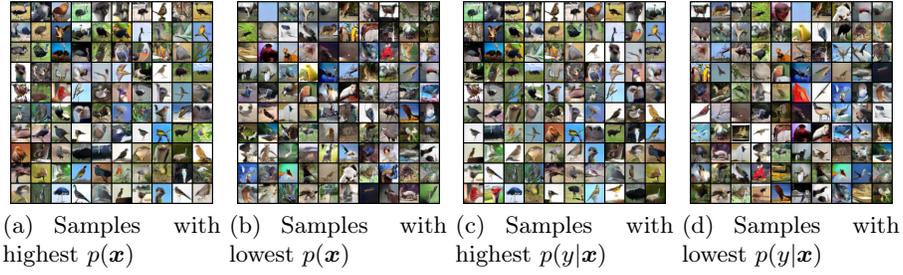


Fig. 11. M-JEM generated class-conditional samples of **Bird**

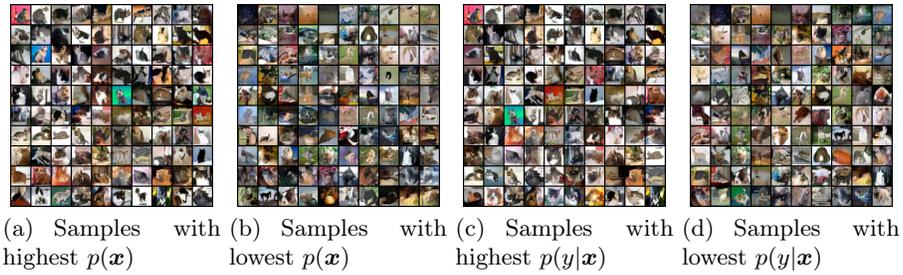


Fig. 12. M-JEM generated class-conditional samples of **Cat**

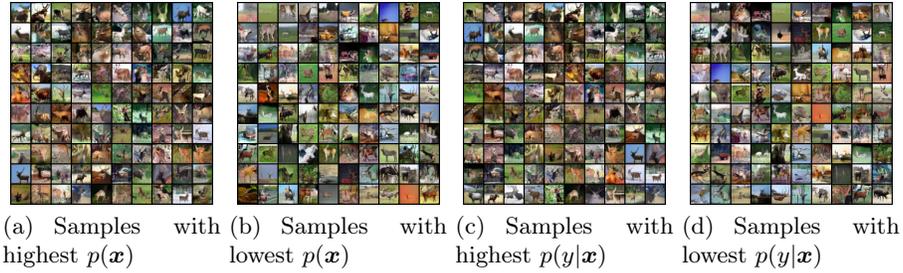


Fig. 13. M-JEM generated class-conditional samples of **Deer**

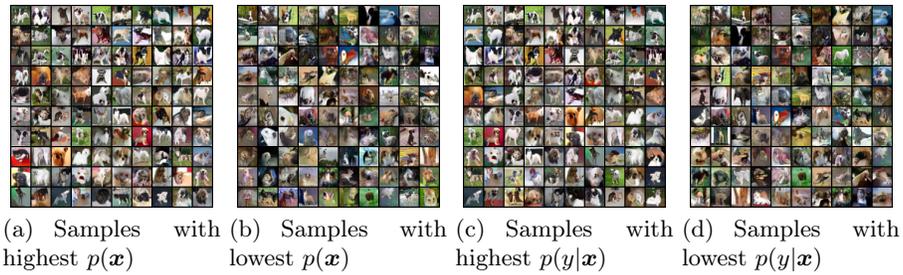
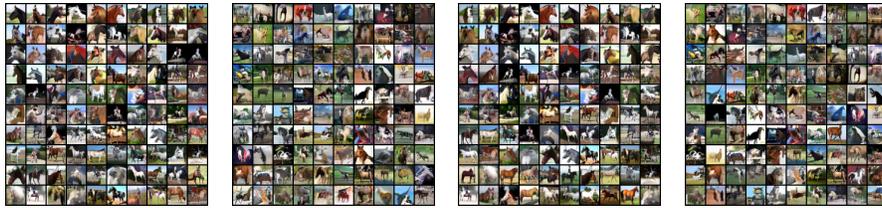


Fig. 14. M-JEM generated class-conditional samples of **Dog**



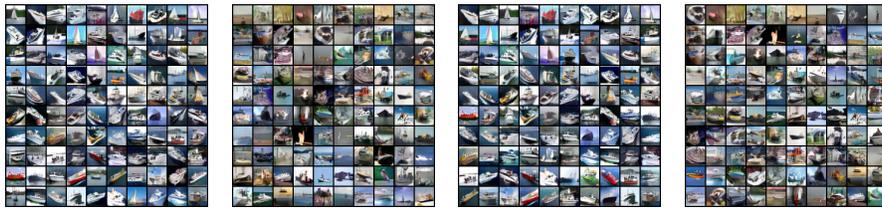
(a) Samples with highest $p(\mathbf{x})$ (b) Samples with lowest $p(\mathbf{x})$ (c) Samples with highest $p(y|\mathbf{x})$ (d) Samples with lowest $p(y|\mathbf{x})$

Fig. 15. M-JEM generated class-conditional samples of **Frog**



(a) Samples with highest $p(\mathbf{x})$ (b) Samples with lowest $p(\mathbf{x})$ (c) Samples with highest $p(y|\mathbf{x})$ (d) Samples with lowest $p(y|\mathbf{x})$

Fig. 16. M-JEM generated class-conditional samples of **Horse**



(a) Samples with highest $p(\mathbf{x})$ (b) Samples with lowest $p(\mathbf{x})$ (c) Samples with highest $p(y|\mathbf{x})$ (d) Samples with lowest $p(y|\mathbf{x})$

Fig. 17. M-JEM generated class-conditional samples of **Ship**