

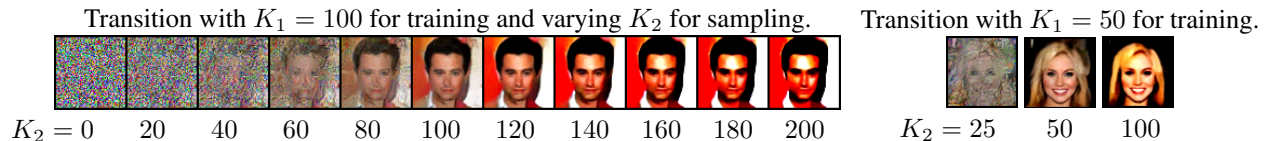
1 **On Learning Non-Convergent Non-Persistent Short-Run MCMC Toward Energy-Based Model.**

2 **Reply to Reviewer 2:** Thank you for the insightful and comprehensive summary of our work.

3 **Q1: About training time. A1:** As you have pointed out, each iteration requires computing  $K$  derivatives of the CNN.  
 4 As an example, 100,000 model parameter updates for  $64 \times 64$  CelebA with  $K = 100$  and  $n_f = 64$  on 4 Titan Xp  
 5 GPUs take 16 hours. We will add such information in revision.

6 **Q2: Dynamic  $K$ . A2:** Following your advice, we conducted experiments with random  $K \in [100, 120]$  for training.  
 7 We can still learn short-run MCMC successfully. This corresponds to residual network with random number of layers.

8 **Q3: Using different  $K$  for training and sampling. A3:** Following your advice, we conducted an experiment on  
 9 CelebA, where we train the model with  $K_1$  and test the trained model by running MCMC with  $K_2$  steps. The Figures  
 10 below depict training with  $K_1 \in \{100, 50\}$  and varied  $K_2$  for sampling. Note, over-saturation occurs for  $K_2 > K_1$ .

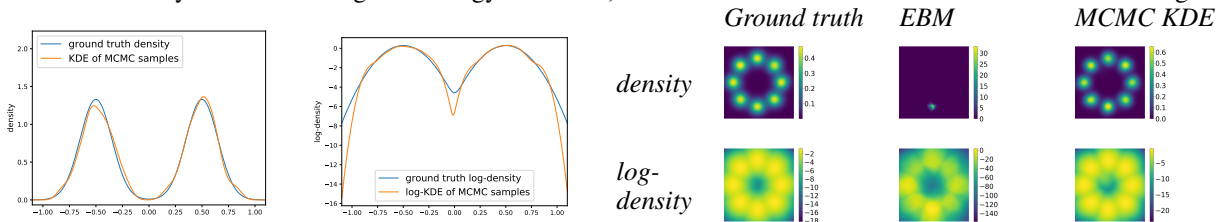


11 **Reply to Reviewer 3:** Thank you for the insightful comments.

12 **Q1: About theoretical justification. A1:** Your comment is well taken, and we shall continue to try our best on  
 13 theoretical understanding. About the generalized moment matching estimator, classical statistical theory shows that the  
 14 estimator is asymptotically unbiased, and the variance of the estimator can also be derived (see Supplementary 2.2).  
 15 The convergence of the learning algorithm follows Robbins-Monro stochastic approximation because  $q_\theta(x)$  can be  
 16 sampled exactly.

17 **Q2: About computing log-likelihood. A2:** Following your advice, we computed log-likelihood. The prior model is  
 18  $z \sim p_0(z)$  which is the uniform distribution. After training, we learn the dynamics  $x = M_\theta(z)$ , where  $M_\theta$  consists  
 19 of  $K$ -steps of gradient descent dynamics (the noise term of the Langevin dynamics is negligible compared to the  
 20 gradient term after learning). Under this flow dynamics  $p(x) = p_0(z)/\det(\partial M_\theta(z)/\partial z)$  (where we used GELU, a  
 21 differentiable version of ReLU). For each observed  $x$ , we can obtain  $z$  as described in Section 3.4. Then we compute  
 22 the Jacobian and its determinant as the product of the eigenvalues of the Jacobian. In our preliminary results, the  
 23 log-likelihood computation is feasible for images of size  $32 \times 32 \times 3$ . For a small batch of 16 images, we have obtained  
 24 a rough, preliminary estimate average 4.12 number of bits per data dimension. We will include log-likelihood results in  
 25 revision. We shall also try to implement Burda et al. (thanks for the reference).

26 **Q3: 1D and 2D toy examples. A3:** Following your advice, we did 1D and 2D experiments. We plot the density and  
 27 log-density of the true model, the learned EBM, and the kernel density estimate (KDE, like histogram) of the MCMC  
 28 samples. The density of the MCMC samples matches the true density closely. The learned energy captures the modes of  
 29 the true density, but is of a much bigger scale, so that the learned EBM density is of much lower entropy or temperature  
 30 (so that the density focuses on the global energy minimum). This is consistent with our theoretical understanding.



31 **Reply to Reviewer 4:** Thank you for the interesting questions.

32 **Q1: About ability to reconstruct. A1:** You are right. It is due to short-run non-mixing, i.e., different starting point  
 33  $z$  leads to different  $x$  after  $K$ -step MCMC, and  $K$  is fixed, so that a mapping is well-defined  $z = M_\theta(x)$ , where  $M_\theta$   
 34 consists of  $K$ -steps of gradient descent dynamics (the noise term of the Langevin dynamics is negligible compared to  
 35 the gradient term after learning).  $M_\theta$  can be considered a flow model.

36 **Q2: Large  $K$  is favored. A2:** We believe larger  $K$  achieves improved synthesis results because larger  $K$  leads to a  
 37 better learned flow model  $M_\theta(z)$  (as a residual network with more layers). Our experience suggests that increasing  $K$   
 38 much further beyond 100 does not lead to much improved results.

39 **Q3: Monte Carlo EM. A3:** Our short-run MCMC always initializes from noise images sampled from the uniform  
 40 distribution. In Monte Carlo EM, one may initialize  $K$ -step MCMC that samples from the posterior distribution of the  
 41 latent variables from the currently generated samples, i.e., running the so-called persistent chains with warm start. Of  
 42 course, one may also initialize from the same noise distribution, i.e., cold start. In that case, the short-run MCMC may  
 43 be interpreted as a variational inference model. We have used this method to learn a latent EBM similar to Boltzmann  
 44 machine but with continuous latent variables, where we use short-run MCMC for both inference and synthesis.