

1 We thank all the reviewers for their insightful comments. We will polish the wording accordingly and will add additional  
2 references as suggested into the next version. We will also release a sample source code example demonstrating the  
3 training and computation of AFVs. Below we address the major concerns, and for comments we do not respond to  
4 explicitly we assume that we agree with the reviewer and will address properly in the revision.

5 1. Training protocol [**Reviewer #1**]: We apologize for the confusion in the training protocol. The simple answer is  
6 that the EBM view does not necessarily change the way the model is trained. As a matter of fact, in the majority of  
7 our experiments, we adopt the standard training procedure of GANs, except for cases where we explicitly test the  
8 effectiveness of the MCMC-inspired objective (see Sec 5.4).

9 A more detailed explanation is as follows. First, we agree that Equation 3 is not a concrete optimization procedure.  
10 However, it does indicate that, in order to train the EBM (D) with the variational trick, G needs to be trained until  
11 convergence to tighten the lower bound on the NLL before updating D. This is in contrast to what is suggested by  
12 Equation 1 (following the convergence analysis of GAN theory), where D needs to be trained until convergence before  
13 updating G. In practice, we approximate both of the max and min optimization problems in Equation 1&3 with a few steps  
14 of mini-batch SGD, as is done in a standard GAN training procedure. As a result, the mini-batch SGD optimization  
15 algorithm can be interpreted as an approximation to either of the two objectives. We will add proper clarifications to the  
16 manuscript.

17 2. Scalability [**Reviewer #2,4**]: Scalability is a practical limitation of our work, because of two reasons. First, computing  
18 Fisher Vectors amounts to taking the gradient of each example w.r.t. all model parameters, which is hard to parallelize  
19 in modern deep learning frameworks. Second, the dimensionality of the induced Fisher Vectors is usually extremely  
20 high, posing a heavy memory demand on GPUs. We made a few initial attempts, one of which is to modify the training  
21 procedure such that we compute only the aggregated Fisher Vector representation for a mini-batch, and accordingly  
22 modify the ground-truth to be the averaged labels of the batch. This resembles an extreme version of mixup, and works  
23 reasonably well when batch size is small, but suffers from performance drop when batch size increases. We suspect that  
24 smarter ways of constructing the batch and averaging the labels might lead to further improvements. For the second fact,  
25 we also tried sampling the parameters to reduce the dimensionality of AFVs, but experienced minor performance drop  
26 in classification accuracy. We believe that the full solution to the scalability issue deserves an independent contribution  
27 and will leave it as future work.

28 3. Loss function [**Reviewer #2,4**]: Our default loss function is least square loss as in LSGAN, with sigmoid activation  
29 on the output of D, except the experiment in Sec 5.4 where we explicitly test the MCMC objective in a new setting. The  
30 reason for such a choice is that it provides the best numerical stability w.r.t. the outputs of D by preventing unnecessary  
31 shifts of D's outputs. We have also tested the hinge loss as done in, e.g., BigGAN, which works equally well w.r.t. the  
32 sampling quality and induced AFV representations, but weakened the smoothness of the Fisher Distance as a monitoring  
33 metric. We will add proper clarifications and discussions.

34 4. The MCMC objective [**Reviewer #2**]: Your understanding is correct: MCMC is never actually performed. We refer  
35 to MCMC because it offers an interpretation of the generator update as approximating one step MCMC. As a result, in  
36 practice we can directly adopt the standard G update rule assuming that each G update is small, mimicking an MCMC  
37 update. In cases where the local update of G is violated, it is useful to explicitly incorporate the proposed MCMC  
38 inspired objective in Equation 7 as a regularizer, as shown in Sec. 5.4. We will make this clear in the paper.

39 5. Additional baselines for the classification experiment [**Reviewer #2,4**]: Per Rev 4's request, during the rebuttal  
40 period we tested the supervised learning performance using the discriminator architecture, by changing the output  
41 dimension of the last layer to 10. With only this change, the supervised learning test accuracy is 86.1%, which is worse  
42 than our AFV + SVM's 89.1%. We then replaced all the Spectral Normalization layers with Batch Normalization and  
43 repeated the experiment, and got a 92.7% accuracy, which exceeds our AFV result. We additionally have conducted the  
44 same experiment on CIFAR100, where AFV+SVM achieves a test accuracy of 67.8%, compared to the supervised  
45 training (with BN) performance 70.3%. Note that 67.8% is also the best result we can find under the pretraining + linear  
46 classifier training setup on CIFAR100. For example, the Deep InfoMax paper reports an accuracy of 49.74%, which is  
47 significantly worse than our result. We will report these experimental results in the paper.

48 6. Approximating Equation 4 [**Reviewer #4**]: The expectation term in Equation 4 is w.r.t. the model distribution  $p_\theta$ . It  
49 is most natural to use the generated samples to approximate it because, according to the EBM view, the generator is  
50 exactly trying to match the model distribution. Empirically, we also found that using the generated examples works  
51 slightly better than using real examples, but the margin is small. We will clarify in the manuscript.