1  We would like to thank all three reviewers for acknowledging our contributions and providing valuable feedback. Please
2  find our responses to your comments below.

3  **Reviewer #1:**

4  Thank you for the positive comments on the novelty of our idea and insightful questions for further improvement.

5  We first characterize the solutions of DEAN. Let $p_x$ and $p_{x'}$ be the distributions of real and fake data; $p_e$ denotes the
6  energy-based distribution. In DEAN, $p_e$ is a bridge connecting $p_x$ and $p_{x'}$. Now we provide two theorems for the
7  characterization. For the IGN, the network is trained to have $p_{x'}$ equal to $p_e$. Please refer to Theorem 1, which is proved
8  based on Theorem 1 of [JXS$^+$17]. For the EGN, $p_e$ is learned to estimate $p_x$. Please see Theorem 2, which is proved
9  according to Theorem 1 of [ZML17] and Theorem 1 of [GPAM$^+$14]. At present, Theorem 2 is proved with $\Lambda(\theta_e)$.
10  Other choices for the energy objective will be left to future works. Detailed proofs of the following theorems will
11  be given in the Supplement of the final version. Different from GANs, which are implicit generative models (IGMs),
12  DEAN can explicitly estimate the underlying distribution of the real data after estimating $\theta_e$ and $\theta_g$.

13  **Theorem 1** *We assume that $\mathcal{D}_{x'}$ is drawn from $p_{x'}$. If the following conditions are satisfied: $\kappa$ is a universal and*
14  *analytic kernel; $\mathbf{E}_{a \sim p_x}\mathbf{E}_{b \sim p_e}\left[s^{\mathrm{T}}(a)s(b)\kappa(a,b) + s^{\mathrm{T}}(b)\nabla_a\kappa(a,b) + s^{\mathrm{T}}(a)\nabla_b\kappa(a,b) + \sum_{i=1}^d \frac{\partial^2 \kappa(a,b)}{\partial a_i \partial b_i}\right] < \infty$ with*
15  *$s(a) = \nabla_a \log p_e(a)$; $\mathbf{E}_{a \sim p_{x'}}\|\nabla_a \log p_e(a) - \nabla_a \log p_{x'}(a)\|^2 < \infty$; $\lim_{\|a\| \to \infty} p_e(a)g(a) = 0$, where $g(\cdot)$ is given*
16  *in Eq. (2) in Section 4.2; for any $J \geq 1$, almost surely $\mathrm{FSSD}[p_e, \mathcal{D}_{x'}] = 0$ if and only if $p_{x'} = p_e$.*

17  **Theorem 2** *Let $\Lambda(\theta_e) = \mathcal{E}(x; \theta_e) + \left[\gamma - \mathcal{E}\left(G(z; \theta_g^*); \theta_e\right)\right]^+$ (please refer to Eq. (1) in Section 4.1 for details). The*
18  *minimum of $\Lambda(\theta_e)$ is achieved if and only if $p_e = p_x$. With the optimized $\theta_e^*$, $\int_{x,z} \Lambda(\theta_e^*)p_x(x)p_z(z)\mathrm{d}x\mathrm{d}z = \gamma$.*

19  Following your suggestion, we compare the powers (successful rejection rates) of MMD, linear-time MMD [GBR$^+$12]
20  and FSSD on toy problems, where MMD is a two-sample test statistic and FSSD is used for the GOF test.
21  We adopt the distributions Gaussian $p(x) = \mathcal{N}(x|0, \mathbf{I}_d)$ and
22  Laplacian $q(x) = \prod_{i=1}^d \mathrm{Laplace}(x_i|0, 1/\sqrt{2})$ for $d = 1, 3$, in
23  which the parameters are set to make $p$ and $q$ have the same mean
24  and variance so that the difference between $p$ and $q$ is subtle. F-
25  SSD shows a higher power to discriminate the subtle difference
26  (Figure 1). For larger sample sizes, the power of MMD is close to
27  that of FSSD. However, in the GAN-type training, the batch size
28  is usually less than 512. As the adversarial training continues, the
29  distribution of the generated data gets closer to the energy-based
30  distribution, and hence the difference becomes subtle. At this
31  time, the power (discriminability) of FSSD for the subtle differ-



Figure 1: Rejection rates for $d = 1$ (left) and $d = 3$.

32  ence becomes important for generating high-quality images. Hopefully we have cleared up your main concerns with
33  these theoretical and experimental discussions. We believe that the DEAN paradigm is promising, being versatile to
34  yield specific training algorithms for different architectures of deep networks in different domains.
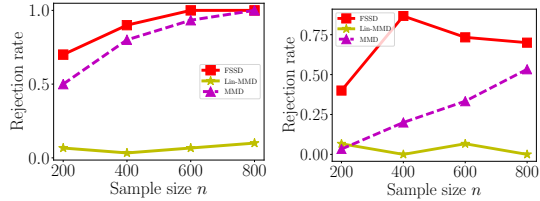
35  **Reviewer #2:**

36  Thank you very much for the encouraging comments and valuable suggestions.

37  Following your recommendation, we will add more discussions in the experimental part to provide takeaways and
38  insights about DEAN. We adopted RBM as the energy function at the initial stage. However, the performance of DEAN
39  with RBM is not comparable to that with autoencoder, so we discarded the results. We will add clarity on this in the
40  final manuscript.

41  **Reviewer #3:**

42  Thank you very much for the positive comments and reasonable doubt.

43  In recent years, there are two emerging families for generative model learning, generative adversarial networks (GANs)
44  and autoencoders (AEs) or variational AEs (VAEs), which are two distinct paradigms and have both received extensive
45  studies. Goodness-of-fit (GOF) tests are a fundamental tool in statistical analysis, dating back to the Kolmogorov test in
46  1933. Our manuscript and [PDB18] both introduce GOF tests into deep generative modeling, but fall into different
47  paradigms: [PDB18] is an AE-based method without adversarial learning while our paper is a GAN-type approach. The
48  HTAE (hypothesis testing AE) in [PDB18] minimized the reconstruction error, but no adversarial learning (min-max
49  adversarial optimization) was involved. The statistic in our manuscript is a kernel-based *nonparametric* GOF statistic.
50  The Shapiro-Wilk test in [PDB18] is a traditional *parametric* GOF statistic for testing normality. Our paper is quite
51  different from [PDB18]. The proposed DEAN with two generators is a pioneering work in the adversarial learning
52  setting. Following your comment, we will cite [PDB18] in the final version.

53  # References

54  [PDB18]  Aaron Palmer, Dipak Dey, and Jinbo Bi. Reforming generative autoencoders via goodness-of-fit hypothesis testing. In
55       *UAI*, pages 1009–1019, 2018.