1 We thank all reviewers for useful feedback.

2 **To Reviewer 1:**

3 *Re. noisy gradient descent (NGD) baseline:* Thanks for the
4 suggestion. We ran this baseline on synthetic (Fig. 1) and image
5 (Fig. 4) data, using sigmoid annealing schedule. We observe NGD to
6 improve over GD on synthetic data. This is intuitive as NGD is better



(a) GD  (b) NGD  (c) AIS-HMC

Figure 1: Error bar plots for GAN after 15k iterations on synthetic data (following Fig. 2 in main paper). AIS-HMC performs best.

7 able to escape some local optima. However, even on synthetic data it performs nowhere close to the AIS-HMC method.
8 *Re. exploration of complex distributions:* Sampling from a multi-modal distribution is challenging. Particularly if the
9 modes are well separated it is important to adequately explore the domain in order not to get stuck in a single mode. To
10 observe this we study the ability to sample from a multi-modal distribution on our synthetic data. We use observation
11 $x_o = x_1 = -1$ which retains an ambiguous $x_2 = 0.5$ or $x_2 = -0.5$. Results are shown in Fig. 3.
12 *Re. other applications:* While other applications are also insightful, we think that synthetic, inpainting and super
13 resolution already display the major challenges: the loss is ragged when optimizing w.r.t. the latent variable. AIS-HMC
14 addresses this ill-posed challenge well, particularly also ambiguity (see Fig. 3 where there exists more than one correct
15 prediction given $x_o = x_1 = -1$).

16 *Re. ablation study:* We perform two studies and show results in Fig. 2: (1) leap



(a)  (b)

Figure 2: Ratio of trials reaching a reconstruction error less than 0.2: (a) leap frog step size over the number leap frog iterations and (b) the number of leap frog iterations over the number of intermediate distributions in AIS.

18 frog step size over leap frog iterations; and (2) leap frog iterations over number
19 of intermediate AIS distributions, i.e., the number of HMC iterations. From (1)
20 we see the method is stable. This is due to HMC adjusting leapfrog step size and
21 acceptance probability. From (2) we note that performance is suboptimal with
22 few intermediate AIS distributions. With one AIS distribution, we run vanilla
23 HMC. This shows that AIS-based HMC has a big advantage over just HMC.
24 *Re. FID:* FID is not an appropriate metric for co-generation. We are interested in
25 retrieving a reconstruction which best fits the given data, hence accuracy matters.
26 Note, diversity does not necessarily exist. In contrast FID assesses diversity (among others). This being said, we do
27 agree that research on how to better assess this task is necessary. This is however beyond the scope.
28 *Re. AIS and necessity:* It is AIS-based HMC because we use an annealing process to move samples from a tractable
29 distribution to the target distribution via a sequence of intermediate steps (line 4 in the algorithm). HMC would
30 directly sample from the target distribution. For complex distributions obtained from GANs we found plain HMC to be
31 challenging (see Fig. 2b when the number of intermediate distributions in AIS=1).
32 *Re. Fig. 4 (paper):* It is first referenced five lines after discussing the experiment (L254). We'll present better. Error
33 bars and additional baselines are shown in Fig. 4 (rebuttal).
34 *Re. SGD:* Thanks for pointing out, it should be GD, we'll revise.
35 *Re. train/test data:* We follow and use prior work, e.g., Progressive GAN. They do not split train/test data. We use 100
36 images for the metric evaluation for all real image tasks. We added error bars to Fig. 4 (rebuttal).

37 **To Reviewer 2:**

38 *Re. optimizing N points:* Thanks for suggesting. This improves slightly compared to gradient descent, but is computa-
39 tionally more expensive. The AIS-HMC method still has a significant edge. See 'MultiOpt+GD' in Fig. 4(c, d).
40 *Re. run-time:* For CelebA data, AIS-HMC takes approximately 13min (0.1min for one HMC step). GD and NGD both
41 take approximately 15min for the 30,000 GD update iterations that we use.
42 *Re. other techniques & clarity:* We'll discuss, add references, e.g., to Dinh et al., and clarify.

43 **To Reviewer 3:**

44 *Re. inpainting & writing:* Please note that this isn't an inpainting paper. We are interested in studying the co-generation
45 task, i.e., how to optimize w.r.t. latent samples. This is more general than inpainting. Agreed, specific methods can be
46 trained for each task, however, this isn't the point. Thanks for pointing out writing, we'll clarify.
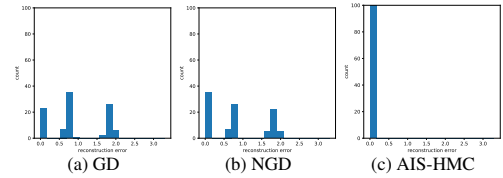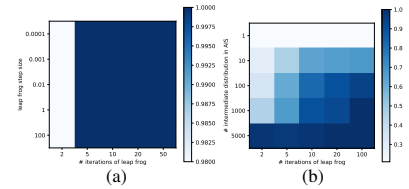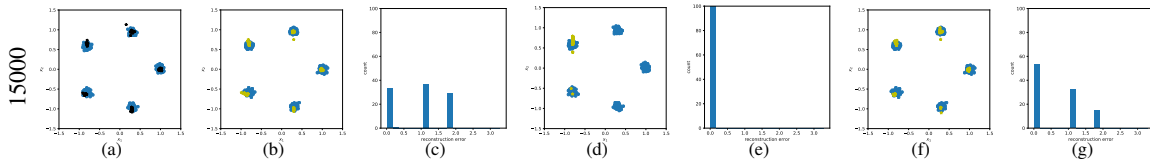


(a)  (b)  (c)  (d)  (e)  (f)  (g)

Figure 3: Columns illustrate: (a) Samples generated with a vanilla GAN (black); (b) GD reconstructions from 100 random initializations; (c) Reconstruction error bar plot for the result in column (b); (d) Reconstructions recovered with Alg. 1; (e) Reconstruction error bar plot for the results in column (d). (f) NGD reconstructions from 100 random initializations; (g) Reconstruction error bar plot for the result in column (f).



(a) higher is better  (b) lower is better  (c) higher is better  (d) lower is better
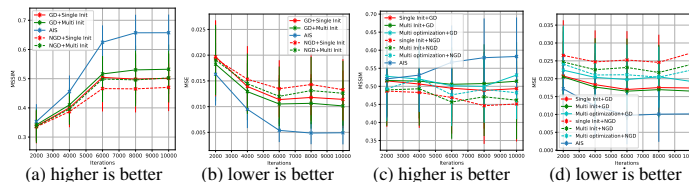
Figure 4: Reconstruction error over Progressive GAN training iterations: (a) MSSIM on CelebA; (b) MSE on CelebA; (c) MSSIM on LSUN; (d) MSE on LSUN.