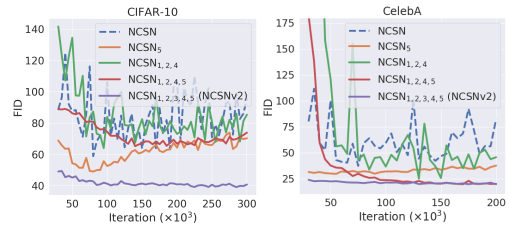


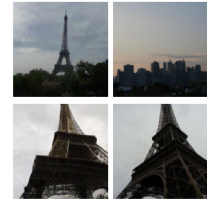
1 We thank all the reviewers for providing valuable feedback in this time of stress. Below we first discuss a new discovery
 2 on evaluation metrics, and then answer specific questions of the reviewers.

3 **A surprising fact of FID computation.** Because of annealed
 4 Langevin dynamics, the samples contain small Gaussian noise that
 5 is imperceptible to human eyes. After our paper submission, we
 6 discovered that this small Gaussian noise—though hard to detect by
 7 humans—can greatly hurt the FID scores. Therefore, we **denoise**
 8 **the samples by running one step of $x \leftarrow x + \sigma_L^2 s_\theta(x, \sigma_L)$ and**
 9 **compute the FID scores again.** We provide the new ablation results
 10 in the right figure, with the extra configuration NCSN_{1,2,4} suggested
 11 by R3. We follow the same checkpoint selection method in Table 5 and provide full FID scores below. **The NCSNv2**
 12 **model now obtains much lower FID scores than NCSN**, which aligns better with our visual inspection of samples.
 13 We are surprised by how the FID scores improve for both NCSN and NCSNv2 though **samples before and after this**
 14 **additional denoising step are the same to naked eyes.** We will include these new results in the revision.



	NCSN (CIFAR-10)	NCSNv2 (CIFAR-10)	NCSN (CelebA)	NCSNv2 (CelebA)
FID	27.44	10.31	17.57	9.69

15 **[R1] Is the model memorizing data (like the Eiffel towers in Figure 1)?** In the paper, we
 16 argued from several perspectives that the model is not memorizing data: (i) The test loss and
 17 training loss are comparable to each other (see Figure 12); (ii) Nearest neighbors in the training
 18 dataset do not look the same as samples from the model (see Section C.4.2); and (iii) The model
 19 can generate samples that smoothly interpolate from one to another (see Figure 7 and Section
 20 C.4.3). In the right figure, we additionally provide nearest neighbors (the right column) in ℓ_2
 21 distance to the two Eiffel towers (the left column) which appeared in Figure 1.



22 **[R1][R2] Whether EMA has a negative impact on performance?** As R1 and R2 noted, EMA stabilizes training but
 23 sometimes may have a slightly worse peak FID score. Because a larger variance gives rise to larger extreme values,
 24 unstable methods naturally lead to a better peak FID score. However, we believe this is an imperfection of the peak FID
 25 metric, rather than an indicator that unstable methods perform better. In fact, as shown in Figure 4 and 11, EMA yields
 26 lower FIDs most of the time and samples with EMA look much more visually appealing than those without EMA.

27 **[R2] Are all techniques needed for scaling to higher resolution?** From the new ablation results above, we observe
 28 that using all techniques leads to the best performance. We agree that for a specific dataset like CelebA it may not be
 29 necessary to use all 5 techniques to get reasonable results, but one key point of our paper is that using all techniques
 30 make the model **work out of the box for a large number of different datasets**, which we demonstrate on many
 31 datasets of different resolutions, including 32^2 , 64^2 , 96^2 , 128^2 and 256^2 .

32 **[R2] When does the RefineNet architecture change?** We hope to clarify one confusion: in the ablation study, only
 33 NCSNv2 uses the new architecture and the others use the old one. The new architecture is necessary for using Technique
 34 3 because it assumes an unconditional score network. We can view the impact of this architecture change by comparing
 35 NCSN_{1,2,4,5} and NCSN_{1,2,3,4,5} (*i.e.*, NCSNv2) in the ablation results.

36 **[R2] Writing issues.** Thanks for pointing them out! We will incorporate your suggestions in the revision.

37 **[R2][R3] Evaluation metrics.** There are many known issues with existing metrics of sample quality, and finding the
 38 right one is still an open problem. We choose FID and HYPE_∞ as an approximation to the real sample quality. The
 39 discrepancy of FIDs in Figure 5 and Table 5 is because FIDs in Table 5 **are the peak FIDs and are computed on**
 40 **more samples.** The HYPE_∞ scores are computed on more than 2000 **uncurated** samples, and are better when closer
 41 to 50. “Fakes Error” is the proportion of fake images perceived as real, and “Reals Error” is the opposite.

42 **[R3] Oversimplified assumptions.** Despite using simplified assumptions, our theory predicts parameters that perform
 43 very well across a large number of complicated real datasets. It is proved by our experiments to be useful and valuable.

44 **[R3] KDE and Technique 1, 2, 4.** When applied to multi-scale KDE as in the
 45 setting of Figure 2, annealed Langevin dynamics will converge to samples that
 46 **exist in the training dataset.** Training an NCSN makes it possible to generalize
 47 to novel samples. We provide the ablation study of Technique 1, 2, 4 in the above
 48 figure (see NCSN_{1,2,4}), showing that they can improve FID over NCSN even
 49 without EMA (Technique 5).

Dataset	Device	Sampling time	Training time
CIFAR-10	2x V100	2 min	22 h
CelebA	4x V100	7 min	29 h
Church	8x V100	17 min	52 h
Bedroom	8x V100	19 min	52 h
Tower	8x V100	19 min	52 h
FFHQ	8x V100	50 min	41 h

50 **[R2][R4] How long the model trains on what hardware, and the sampling speed.** We provide the statistics in the
 51 above table, and will add it to the paper in the next revision. The sampling time is for one mini-batch.