

1 Thank you for your comments. Here we address the most important ones (we can and will also fix the others):

2 **Rev#1/Rev#4: What’s the difference between Deeploss-VGG/-Squeeze and the loss proposed in [29] (LPIPS)?** It is the
3 same metric, as noted in L140. We wanted a consistent naming scheme in the paper, but see that this can be confusing.
4 We consider renaming it to LPIPS-VGG and LPIPS-Squeeze.

5 **Rev#1: Comparison to MSE for VAE** In our experience there are no big differences on MNIST between L_2 and MSE
6 and the SSIM shown (which is in general considered the superior loss for images). If any, the generated images look
7 slightly less crisp, however this effect is often dominated by the quality of hyper-parameter tuning.

8 **Quantitative analysis on other score-functions** We will provide a measure in the updated paper. The choice is a bit
9 difficult because, for example, an LPIPS based measure will work best on the models trained with a variant of LPIPS.

10 **Rev#2: Evidence that loss actually guides generative models towards better image generations** The proposed metric
11 does lead to *much* better results, see Section 4.2, figures 4, 5, and in the supplement Section E, figures E.8–E.11.

12 **Correctness of results in relation to “Deep Feature Consistent Variational Autoencoder”** In the article by Hou et al. the
13 weights are not trained but fixed to $\omega_{l_c} = 1/C_l$ in our notation. When the weights are adapted based on another dataset,
14 the resulting losses will be different. This is especially true in our case. While the dataset used for tuning is relatively
15 large, it still covers only a small subset of relevant transformations and the number of tunable parameters is large. Thus,
16 overfitting to the dataset can introduce, e.g., warping artefacts.

17 **Rev#4: Comparison to “A General and Adaptive Robust Loss Function”, Jonathan T. Barron, CVPR, 2019.** Thank you, we were not
18 aware of that interesting and relevant publication. There are many
19 difference to our approach: Barron performs a 2D DCT over the
20 entire image while we use a blockwise FFT, which allows us to also
21 consider phase differences. This improves perceptual accuracy signi-
22 ficantly compared to DCT. We weight the DCT/FFT frequencies
23 resulting in <200 weights. Barron learns a ‘robustness’ value for each
24 separate DCT frequency. Due to not using block-DCT, this results in
25 a lot of parameters (e.g., 49152 trainable parameters on $128 \times 128 \times 3$
26 images). We also use the YUV/YCbCr color space, but we weight the color channels by learning the importance of each,
27 while Barron weights them equally. Both approaches learn a “robustness” parameter determining the significance of
28 outliers (α vs. p). Barron learns this parameter during training on the generative task, but has to add some regularization
29 to make this work. We learn the perceptual parameters on a perceptual dataset, independent of the generative task
30 and no regularization is necessary. We conducted experiments using the CVPR 2019 method with the code thankfully
31 provided by the author, e.g., see figures R.1 and R.2 (we will add more to the paper, also showing generated images).
32 *The method performs well, but clearly worse than Watson-DFT.*

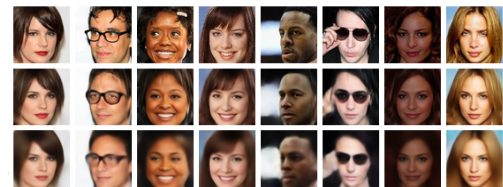


Figure R.1: Reconstructions using VAE. 1st row: Ground truth; 2nd: Watson-DFT, 3rd: CVPR 2019 method (please zoom in).

34 **Performance differences in Figure-6 & user-study** *This is an impor-*
35 *tant aspect/finding of our study, please revisit L240–L270.* Consider
36 what task is solved and measured in Figure 6: a dataset is generated
37 by applying a certain set of transformations to images. The test-set
38 is not an unbiased estimate of performance in a real application as it
39 does not consider all relevant transformations. Moreover, the gener-
40 ative task differs from the test-set insofar as the VAE, similar to the
41 adversary in a GAN, tends to find the weaknesses in the loss-function
42 in order to maximize the similarity of $q(z|x)$ to $N(0, I)$. A user-study
43 would be nice to have, but may not provide insights into the models.
44 We believe that the differences between the models are so large that
45 there is no need for a user-study to decide which images look better –
46 we assume that all reviewers agree on the obvious visual differences (on
47 random samples, more can be generated using the software provided).

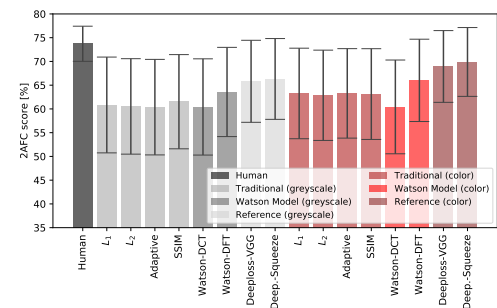


Figure R.2: Updated Figure 6 of the submission, now including L_1 and Watson-DCT and the CVPR 2019 method (“Adaptive”).

48 **Free parameters & Ablation study of model components** A full list of free parameters is given in L137. The maximum
49 number of parameters of our method is 135. Earlier during development we compared DFT and DCT within our model.
50 With DCT the model performed significantly worse in all tasks, on par with the other models. We included this result in
51 Fig. R.2 as “Watson-DCT”.

52 **MCMC and the loss as probability distribution** We did not want to claim that one can use MCMC if the model is not a
53 valid probability distribution. We will clarify this part. When the loss is a valid unnormalized log-probability, we can
54 use standard MCMC techniques like HMC to sample from $p(x|z)$.