1 **Author Response for "Generating Diverse High-Fidelity Images with VQVAE-2"**
2 We thank the reviewers for the detailed and constructive feedback. All reviewers were impressed by the quality of our
3 samples. R2 had positive remarks about the significance of our method and the thoroughness of our evaluation. R3 was
4 satisfied with the clarity of the writing. There were however some concerns about additional results, architecture details
5 and novelty of the approach. We hope that we have addressed the requested clarifications below, explaining how each
6 will be improved in the final paper. We believe these clarifications will resolve all reviewers' concerns, but would be
7 happy to consider any additional suggestions.
8 **General response regarding to:**
9 **R1 -** *Name of the method*: The name VQ-VAE was coined by previous authors, so we decided to adopt the same name,
10 rather than using a new one. Moreover, it is possible to frame VQ-VAE in a variational framework as well (with a delta
11 posterior and a uniform prior), which is discussed in the original VQ-VAE paper.
12 **R1 -** *h_bottom encoding everything*: **A** h_bottom is preceded by an Information-Bottleneck (similar to the KL of a
13 regular VAE) so every latent can at most encode $\log_2(N)$ bits, where N is the size of the VQ codebook. In our case, this
14 is 9bits per latent. There is one latent for every 16 pixels, each having 3 color channels with 8 bits each, so this yields
15 a compression factor of 4*4*3*8/9=42.67. **B** Adding h_top results in much better reconstruction MSE and sharper
16 reconstructions. We also visualize reconstruction from h_top only (Fig. 3 in the paper and Fig. 5 in the appendix),
17 showing that h_top has indeed encoded quite a lot of information.
18 **R1 -** *Model remembering data*: The VQ-VAE is able to reconstruct test set images equally well as training images
19 (MSE, and qualitatively), which demonstrates that it generalizes to unseen data. Similarly, for the PixelCNN, the NLLs
20 of test data are comparable to those of training data.
21 **R1 -** *Interpolations*: There is indeed no simple way to do interpolations, which We will clarify in the final version.
22 **R1 -** *Speed*: Generation is slower than GANs, but faster than other autoregressive approaches that model images in the
23 pixel space (about 45x faster). We also implemented incremental sampling (as in Paine et al. arxiv.org/abs/1611.09482)
24 to cache intermediate activations that can considerably reduce sampling time. We will add a comparison in the final
25 version.
26 **R1 -** *Objective with stop-gradient not elegant*: As noted in the paper, in equation 3, we use the Exponential Moving
27 Average version which does not use stop-gradients. The loss in equation 2 is included for sake of completeness. These
28 are both different neural implementations of the K-means algorithm, which has a long established track record in many
29 areas of machine learning. We would also like to point out that elegance is a subjective matter, and simplicity is a form
30 of elegance we strove for in this work: indeed, VQ-VAE is quite simple and can be implemented in just a few lines of
31 code.
32 **R2 -** *Detailed architecture*: We agree that the architecture description could be more detailed. We will make sure that
33 our architecture is thoroughly specified in our final version and we will include all details and hyperparameters.
34 **R2 -** We will fix minor details, also cite [A], [B] to emphasize the connection with lossy-compression.
35 **R3 -** *Novelty / This paper is not the first to achieve that*: We are not aware of any prior works that show comparable
36 sample quality to BigGAN (which is SOTA) while having better diversity in any model class (let alone among
37 likelihood-based methods). For faces, the best prior works (ie, Glow and SPN) used a much simpler dataset CelebaHQ
38 with 256x256 resolution. Still their their samples look less realistic and have lower fidelity than our 1024x1024 samples
39 from more complex FFHQ. The only other model to achieve this has been StyleGAN, which also has the discussed
40 diversity limitations of GANs.
41 **R3 -** *Ablations wrt. model size*: We will add ablations of our model wrt. model size and batch size, but the result is that
42 larger models get better results. Comparison with other works, however, shows that scaling up is necessary but not
43 sufficient: our model gets better results compared to models with similar (or larger) size, batch-size, and compute
44 requirements: Menick et al. 2018, Defauw et al. 2019. The same applies for BigGAN.
45 **R3 -** *No quantization*: The model does not work at all without quantization. We will add this ablation in the Appendix
46 (with a reference from the main text).
47 **R3 -** *BPD measurements*: As R2 has noted, this model is inspired by lossy-compression where performance is usually
48 characterized with rate-distortion curves. We apply log-likelihood based methods in a compressed lossy space, thus not
49 having to model imperceptible details in images. This is where benefits like speed, global coherence, etc., come from.
50 We do report BPD in the latent spaces. Trying to go back to the pixel-domain would defeat the main purpose of the
51 method. That said, if there is truly interest in this metric, it is straightforward to estimate and add it to the final version
52 of the paper.
53 **R3 -** *Classifier based rejection sampling*: All samples in the paper are without the rejection sampling except for Fig 8
54 and 9 in the appendix where we aim to illustrate the effect of various rejection thresholds. Similarly the numbers in
55 Table 1 do not use this. CBRS is only used for the P/R and FID/IS curves in Fig. 5. of the main text.
56 **R3 -** *Nearest Neighbours*: As noted in the paper and in our response to R1, our model can be directly assessed for
57 overfitting by comparing train and test NLL. Nevertheless, we will include nearest neighbours in the pixel and VGG
58 spaces in the final version.
59