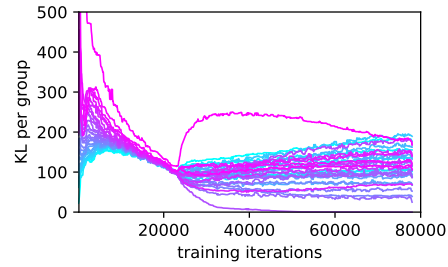(a) Reconstruction results (best seen when zoomed in).



(b) Average KL per group.

Figure 1: (a) Input on the left and reconstructed image on the right for CelebA HQ 256. (b) KL per group on CIFAR10.

1  We would like to thank all the reviewers for positive and constructive feedback.

2  **R1** ———————————————————————————————————————————————

3  **Reconstruction:** The reconstructed images in NVAE are indistinguishable from the training images (see Fig. 1(a)).

4  **Discussing additional SOTA hybrid models:** Thanks for pointing this out. We will include this in our final version

5  **Non-dimension-reduction and non-disentangling design:** Since our goal is to maximize (a lower bound on) the
6  marginal data log-likelihood, hierarchical dependencies help us build expressive approximate posteriors, which often
7  result in better generative performance. However, when the training goal shifts towards disentangled representation
8  learning, we can sacrifice generative performance for representation learning (See $\beta$-VAE, "Fixing a Broken ELBO").

9  **Non-face datasets such as ImageNet:** We haven't explored the ImageNet dataset beyond the 32x32 version. However,
10  our hypothesis is that ImageNet may correspond to a high-entropy distribution (compared to face images) which may
11  require even bigger VAE models. GANs are perhaps less prone to this, as they may drop modes without being penalized.

12  **R2** ———————————————————————————————————————————————

13  **Is the data conditioned on all $z$'s:** Yes, because the representation at the bottom of the top-down model (Fig. 2(b)) is
14  a function of all $z$'s.

15  **Parametrization for $\sigma$:** We parameterize $\sigma$'s in their log space, and we limit $\log \sigma$ to be in [-5, 5] which is much larger.

16  **Training curves:** Fig. 10 in the supplementary material demonstrates training stability with spectral regularization.

17  **Posterior collapse:** Since we are using more latent variables than the data dimensionality, it is natural for VAE objective
18  to turn off many latent variables. However, our KL balancing mechanism (Sec. A in the appendix) stops the hierarchical
19  groups from turning off. In Fig. 1(b), you can see KL per group in CIFAR10 (for 30 groups). Note how most groups
20  obtain a similar KL on average, and only one group is turned off. We apply KL balancing mechanism only during KL
21  warm-up (the first $\sim$ 25000 iterations). In the remaining, we are using ELBO without any KL balancing (Eq. 1).

22  **Importance weighted (IW) AE:** We <u>do not</u> train our model with the IW bound. The comparison to the previous work
23  is fair in this regard as we also use ELBO for training (see Eq. 1).

24  **Autoregressive decoder + NVAE:** Autoregressive decoders 1) are often slow to sample from and 2) tend to turn off
25  latent variables. A careful study is required in this space.

26  **Sampling time compared to pixelCNN++:** On CIFAR10 using a Titan V GPU and batch size of 16, NVAE takes 5.1
27  ms per image vs. 8,898.75 ms per image for pixelCNN++. NVAE is $\sim$1700x faster than pixelCNN++.

28  **R3** ———————————————————————————————————————————————

29  **Reproducibility:** Our current top priority is to release the code publicly with clear instructions to reproduce the results.

30  **Test likelihood after BN adjustment:** Thanks for suggesting this. Since the generative model changes after re-
31  adjusting the BN layers, this requires encoder retraining. We'll add a careful examination of this to the final version.

32  **R4** ———————————————————————————————————————————————

33  **Lack of novelty:** NVAE offers expressive networks and it enables stable training of deep hierarchical VAEs for large
34  images. We believe NVAE has the potential of becoming the backbone of future VAEs with complex statistical models.

35  **Efficacy of depthwise conv:** Depthwise conv has less computational complexity, however, it is still slow on current
36  hardware because of memory bottleneck. We agree with you that depthwise conv seems less efficient, but in our early
37  experiments, the generative quality of regular conv could not match it, even with more channels or longer training.

38  **VQ-VAE:** Thanks for pointing this out. We will update our description of VQ-VAE to indicate that it was originally
39  motivated by deriving a lower bound on log-likelihood.