
Distribution Matching in Variational Inference

Mihaela Rosca Balaji Lakshminarayanan Shakir Mohamed
{mihaelacr,balajiln,shakir}@google.com

DeepMind

Abstract

With the increasingly widespread deployment of generative models, there is a mounting need for a deeper understanding of their behaviors and limitations. In this paper, we expose the limitations of Variational Autoencoders (VAEs), which consistently fail to learn marginal distributions in both latent and visible spaces. We show this to be a consequence of learning by matching conditional distributions, and the limitations of explicit model and posterior distributions. It is popular to consider Generative Adversarial Networks (GANs) as a means of overcoming these limitations, leading to hybrids of VAEs and GANs. We perform a large-scale evaluation of several VAE-GAN hybrids and analyze the implications of class probability estimation for learning distributions. While promising, we conclude that at present, VAE-GAN hybrids have limited applicability: they are harder to scale, evaluate, and use for inference compared to VAEs; and they do not improve over the generation quality of GANs.

1 INTRODUCTION

This paper focuses on the challenges of training Variational Autoencoders (VAEs) (Kingma and Welling, 2013, Rezende et al., 2014): the struggles of matching the marginal latent posterior distribution with the prior, the difficulties faced in explicitly specifying latent posteriors, and the lack of likelihoods that capture the semantic similarity of data. To overcome the limitations of VAEs, it has become natural to consider borrowing the strengths of another popular type of generative algorithm, Generative Adversarial networks (GANs) (Goodfellow et al., 2014), resulting in a fusion of VAEs and GANs (Larsen et al., 2016, Makhzani et al., 2015, Mescheder et al., 2017, Pu et al., 2017, Rosca et al., 2017, Srivastava et al., 2017). What are

the limitations of VAEs at present? How can GANs help? Do VAE-GAN hybrids address the limitations currently being experienced? The aim of this paper is to offer an answer to these questions.

Generative models currently have a wide range of applications in dimensionality reduction, denoising, reinforcement learning, few-shot learning, data-simulation and emulation, semi-supervised learning and in-painting (Higgins et al., 2017b, Kingma et al., 2014, Mathieu et al., 2016, Nguyen et al., 2016, Rezende et al., 2016, Salimans et al., 2017, Smith, 2002), and they continue to be deployed in new domains, from drug-discovery to high-energy physics. It thus becomes vital for us to explore their strengths and shortcomings, and deepen our understanding of the performance and behavior of generative models.

We will contrast conditional distribution matching in VAEs with marginal distribution matching in VAE-GAN hybrids. We compare the performance of explicit distributions in VAEs with the use of implicit distributions in VAE-GAN hybrids. And we measure the effect of distributional assumptions on what is learned. We will make the following contributions:

- We show that VAEs fail to match marginal distributions in both latent and visible space and that powerful explicit model distributions and powerful explicit posteriors do not improve marginal distribution matching in VAEs. This is prevalent across datasets, models and latent dimensionality.
- We systematically evaluate existing and new VAE-GAN hybrids as promising avenues to improve variational inference, but show that since they use classifier probabilities to estimate density ratios and learn implicit distributions, they lack an accurate estimate for the likelihood bound that can be used for model evaluation and comparison.
- We uncover the effect of marginal distribution matching in latent space on latent representations and learned posterior distributions, and show that VAE-GAN hybrids are harder to scale to higher latent dimensions than VAEs.

These contributions will lead us to conclude that while VAE-GAN hybrids allow for posterior inference in implicit generative models, *at present*, VAE-GAN hybrids have limited applicability: compared to VAEs, their use of classifier probabilities makes them harder to scale, evaluate, and use for inference; compared to GANs, they do not improve sample generation quality.

2 CHALLENGES WITH VAES

One goal of generative models is to match the unknown distribution of data $p^*(\mathbf{x})$ to the distribution learned by a model. In latent variables models, $p_{\theta}(\mathbf{x})$ is defined via a latent variable \mathbf{z} and the hierarchical model:

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}. \quad (1)$$

The integral in Equation 1 for computing $p_{\theta}(\mathbf{x})$, in general, intractable, making it hard maximize the marginal likelihood of the model under the data, $\mathbb{E}_{p^*(\mathbf{x})}[\log p_{\theta}(\mathbf{x})]$. To overcome this intractability, variational inference introduces a tractable lower bound on $p_{\theta}(\mathbf{x})$ via a variational distribution $q_{\eta}(\mathbf{z}|\mathbf{x})$:

$$\begin{aligned} & \log p_{\theta}(\mathbf{x}) - \text{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] \quad (2) \\ \log p_{\theta}(\mathbf{x}) &\geq \underbrace{\mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{likelihood term}} - \underbrace{\text{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]}_{\text{KL penalty}} \quad (3) \end{aligned}$$

The training objective is to learn the model parameters θ of the distribution $p_{\theta}(\mathbf{x}|\mathbf{z})$ and the variational parameters η of $q_{\eta}(\mathbf{z}|\mathbf{x})$ to maximize the evidence lower bound (ELBO):

$$\mathbb{E}_{p^*(\mathbf{x})} [\mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]] \quad (4)$$

Distributions $q_{\eta}(\mathbf{z}|\mathbf{x})$ and $p_{\theta}(\mathbf{x}|\mathbf{z})$ can be parametrized using neural networks and trained jointly using stochastic gradient descent as in Variational Autoencoders (Kingma and Welling, 2013, Rezende et al., 2014). In this setting, we refer to $q_{\eta}(\mathbf{z}|\mathbf{x})$ as the encoder distribution, and to $p_{\theta}(\mathbf{x}|\mathbf{z})$ as the decoder distribution.

The ideal learning scenario does *marginal distribution matching*, in visible space: matching $p_{\theta}(\mathbf{x})$ to the true data distribution $p^*(\mathbf{x})$. Variational inference, using a lower bound approximation, achieves this indirectly using *conditional distribution matching* in latent space: minimizing $\text{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})]$ in Equation 2 and searching for a tight variational bound.

Marginal distribution matching in *visible* space, and conditional distribution matching in *latent* space, are tightly related through marginal distribution matching in latent space. By doing a ‘‘surgery on the ELBO’’, Hoffman and Johnson (2016) showed how variational

inference methods minimize $\text{KL}[q_{\eta}(\mathbf{z})||p(\mathbf{z})]$ when maximizing the ELBO, since:

$$\begin{aligned} & \mathbb{E}_{p^*(\mathbf{x})}\text{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] = \quad (5) \\ & \text{KL}[q_{\eta}(\mathbf{z})||p(\mathbf{z})] + \int q_{\eta}(\mathbf{z}|\mathbf{x})p^*(\mathbf{x}) \log \frac{q_{\eta}(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} d\mathbf{z}d\mathbf{x} \end{aligned}$$

The integral in (5) is a mutual information and is non-negative. Minimizing $\mathbb{E}_{p^*(\mathbf{x})}\text{KL}[q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$ therefore minimizes a lower bound on $\text{KL}[q_{\eta}(\mathbf{z})||p(\mathbf{z})]$, i.e. matching the marginal distributions.

There are two reasons to be interested in marginal distribution matching. Firstly, we can unveil a failure to match $q_{\eta}(\mathbf{z})$ to $p(\mathbf{z})$, and $p_{\theta}(\mathbf{x})$ to $p^*(\mathbf{x})$ in VAEs. This failure can be associated to the use of conditional distribution matching and explicit distributions. Our goal will be to explore the solutions to overcome these failures provided by marginal distribution matching and implicit models. Secondly, the marginal divergence $\text{KL}[q_{\eta}(\mathbf{z})||p(\mathbf{z})]$ can be used as a metric of VAE performance since it captures simultaneously two important characteristics of performance: the ability of the learned posterior distribution to match the true posterior, and the ability of the model to learn the data distribution.

An alternative examination of VAEs as a rate-distortion analysis, achieved by using a coefficient in front of the KL term in the ELBO (e.g., Higgins et al. (2017a)), is the closest to our work. Alemi et al. (2017) show that simple decoders have a smaller KL term but have high reconstruction error, whereas complex decoders are better suited for reconstruction. Their investigation does not consider implicit distributions and alternative approaches to distribution matching, and is complementary to this effort.

3 QUANTIFYING VAE CHALLENGES

To motivate our later exploration of marginal distribution matching and implicit distribution in variational inference, we first unpack the effects of conditional distribution matching and explicit distributions. We show that VAEs are unable to match the marginal latent posterior $q_{\eta}(\mathbf{z})$ to the prior $p(\mathbf{z})$. This will result in a failure to learn the data distribution, and manifests in a discrepancy in quality between samples and reconstructions from the model. To the best of our knowledge, no other extensive study has been performed showing the prevalence of this issue across data sets, large latent sizes, posterior and visible distributions, or leveraged this knowledge to generate low posterior-probability VAE samples. Hoffman and Johnson (2016) initially showed VAEs with Gaussian posterior distributions,

Bernoulli visible distribution and a small number of latents trained on binary MNIST do not match $q_\eta(\mathbf{z})$ and $p(\mathbf{z})$.

We train VAEs with different posterior and model distributions on ColorMNIST (Metz et al., 2017), CelebA (Liu et al., 2015) at 64x64 image resolution, and CIFAR-10 (Krizhevsky, 2009). Throughout this section, we estimate $q_\eta(\mathbf{z})$ via Monte Carlo using $\frac{1}{N} \sum_{n=1}^N q(\mathbf{z}|\mathbf{x}_n)$ (pseudocode in Appendix E).

3.1 Effect of latent posterior distribution

Powerful explicit posteriors given by Real Non Volume Preserving (RNVP) normalizing flows (Dinh et al., 2016) do not improve marginal distribution matching in VAEs over simple diagonal Gaussian posteriors. This result, shown in Figure 1 which reports $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ for trained VAEs, is surprising given that RNVP posteriors are universal distribution approximators, and have the capacity to fit complex posteriors $p_\theta(\mathbf{z}|\mathbf{x})$ and perhaps suggests that posterior distributions are not the biggest bottleneck in VAE training.

3.2 Effect of the visible distribution

The choice of visible distribution $p_\theta(\mathbf{x}|\mathbf{z})$ also affects distribution matching in latent space. We trained the same VAE architecture using Bernoulli and Quantized Normal visible distributions on CIFAR-10. A Quantized Normal distribution is a normal distribution with uniform noise $\mathbf{u} \in [0, 1]$ added to discrete pixel data (Theis et al., 2016). Since a Bernoulli distribution provides the same learning signal as a QuantizedNormal with unit variance (formal justification in Appendix G), the QuantizedNormal distribution has a higher modeling capacity, leading to better data reconstructions. However, this does not result an increased sample quality: better reconstructions force $q_\eta(\mathbf{z}|\mathbf{x})$ to encode more information about the data, making the KL between $q_\eta(\mathbf{z}|\mathbf{x})$ and the data agnostic $p(\mathbf{z})$ higher, thus increasing the gap between sample and reconstruction quality. Apart from visually assessing the results (see Appendix G), we can now quantify the different behavior using $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$, which is 44.7 for a Bernoulli model, and 256.7 for the same model with a QuantizedNormal distribution. We saw similar results when comparing Categorical and Bernoulli distributions.

3.3 Low probability posterior samples

We exploit the failure of VAEs to match marginal distributions in latent space to purposefully generate ‘bad’ samples: samples obtained from prior samples z that have low probability under the marginal posterior $q_\eta(\mathbf{z})$. Such bad samples, shown in Figure 2, exhibit common

characteristics such as having thickened-lines for ColorMNIST and strong white backgrounds for CelebA and CIFAR-10. We can look at *where* in latent space the samples with low probability under the marginal are. Figure 3 (right) is a t-SNE visualisation which shows that the latents which generate the low posterior VAE samples in Figure 2 are scattered throughout the space of the prior and not isolated to any particular regions, with large pockets of prior space which are not covered by $q_\eta(\mathbf{z})$.

Since these ‘bad’ samples are unlikely under the true distribution, we expect a well trained model to distinguish easily, using the likelihood bound, typical images (which come from the data) from atypical samples (generated from regions of the marginal distribution with low-probability). Figure 3 (left) shows that on ColorMNIST the model recognizes that the low posterior samples and their nearest neighbors have a lower likelihood than the data, as expected. But for CIFAR-10 and CelebA the model thinks the low posterior samples are more likely than the data, despite the samples being quantitatively and qualitatively different (see Appendix A for more samples and metrics).

The latent space spread of low posterior samples, together with the inability of VAEs to learn a likelihood that reflects the data distribution, show a systematic failure of these models to match distributions in latent and visible space.

4 GANS IN LATENT SPACE

If it is conditional divergence minimization and explicit posteriors that result in a failure of VAEs to match marginal distributions, then we are strongly motivated to explore other approaches. We turn to another popular type of algorithm, generative adversarial networks (GANs), and explore *the density ratio trick*, which they use as a tool for marginal distribution matching and implicit distributions in variational inference, leading to VAE-GAN hybrids.

4.1 Distribution matching with density ratios

The density ratio trick (Goodfellow et al., 2014, Goodfellow, 2014, Sugiyama et al., 2012) leverages the power of classifiers in order to estimate density ratios. Under the assumption that we can train a perfect binary classifier \mathcal{D} to associate the label $y = 1$ to samples from $p_1(\mathbf{x})$ and the label $y = 0$ to samples from $p_0(\mathbf{x})$, the following holds:

$$\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} = \frac{p(y=1|\mathbf{x})}{p(y=0|\mathbf{x})} = \frac{\mathcal{D}}{1-\mathcal{D}} \quad (6)$$

Conveniently, this approach only requires distribution samples, without the explicit forms of the two

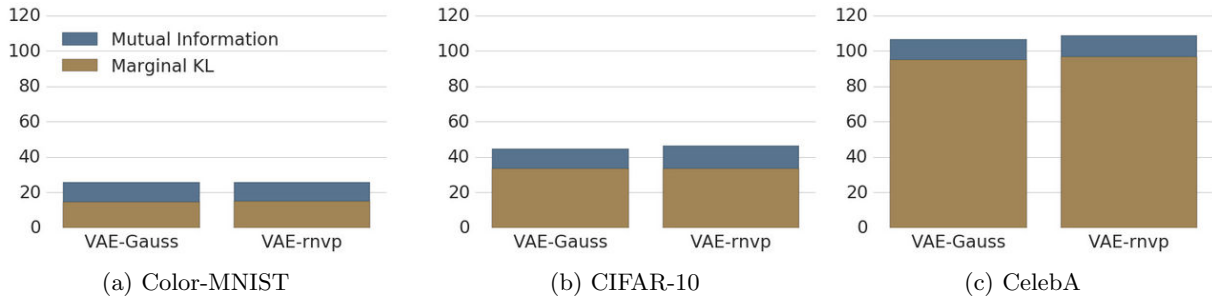


Figure 1: Marginal KL and mutual information of VAEs with Bernoulli visible distribution and Gaussian and RNVP posteriors, adding up to the KL term in the ELBO. $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ reflects the complexity of the data: all models use 160 latents, but complex datasets exhibit higher values.

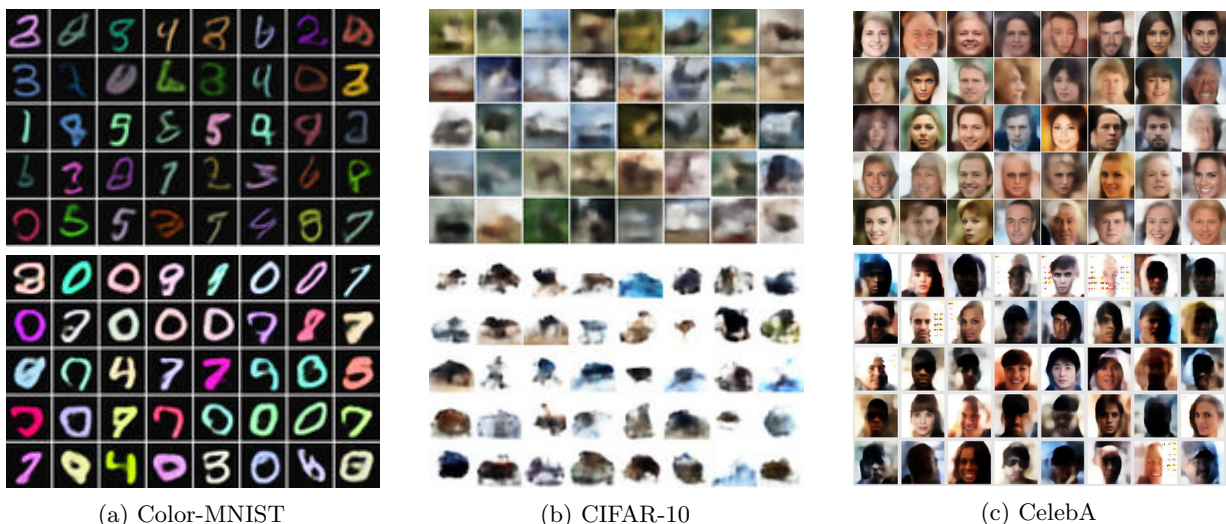


Figure 2: VAE samples (top) and low posterior VAE samples (bottom). More samples in Appendix A.

distributions. GANs (Goodfellow et al., 2014) use the density ratio trick to directly match the data distribution $p^*(\mathbf{x})$ with the marginal model distribution $p_\theta(\mathbf{x})$. By making use of *implicit latent variable models* (Mohamed and Lakshminarayanan, 2016) GANs do not require observation likelihoods - they define $p_\theta(\mathbf{x})$ by specifying a deterministic mapping $G : \mathbf{z} \rightarrow \mathbf{x}$ used to generate model samples. GANs learn via an adversarial game, using a discriminator to distinguish between generated samples and real data. In the original formulation, the training was given by the min-max bi-level optimization with value function¹:

$$\min_G \max_D \mathbb{E}_{p^*(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{p(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

When the density ratio trick is used for learning - like in GANs - the ratio estimator cannot be trained to optimality for each intermediate model distribution, since this would be computationally prohibitive. In

¹In all our experiments we use the alternative non-saturating generator loss $\mathbb{E}_{p_\theta(\mathbf{x})} [-\log \mathcal{D}_\phi(\mathbf{x})]$ (Goodfellow et al., 2014).

practice, the model and the estimator are trained jointly using alternating gradient descent (Fedus et al., 2018, Goodfellow et al., 2014). Since KL divergences are an expectation of density ratios, the density ratio trick opens the door to implicit posteriors and marginal distribution matching in latent space for variational inference.

Implicit variational posteriors. The choice of posterior distribution in VAEs is limited by the use of the choice of $\log q_\eta(\mathbf{z}|\mathbf{x}_n)$ in $\mathbb{E}_{p^*(\mathbf{x})} [\text{KL}[q_\eta(\mathbf{z}|\mathbf{x}_n)||p(\mathbf{z})]]$. Since the density ratio trick only requires samples to estimate $\text{KL}[q_\eta(\mathbf{z}|\mathbf{x}_n)||p(\mathbf{z})]$, it avoids restrictions on $q_\eta(\mathbf{z}|\mathbf{x})$ and opens the door to implicit posteriors. However, replacing each KL divergence in equation 4 with a density ratio estimator is not feasible, as this requires an estimator per data point. AdversarialVB (Mescheder et al., 2017) solves this issue by training one discriminator to estimate all ratios $q_\eta(\mathbf{z}|\mathbf{x})/p(\mathbf{z})$ for $x \sim p^*(\mathbf{x})$ by distinguishing between pairs $(\mathbf{x} \sim p^*(\mathbf{x}), \mathbf{z} \sim p(\mathbf{z}))$ and $(\mathbf{x} \sim p^*(\mathbf{x}), \mathbf{z} \sim q_\eta(\mathbf{z}|\mathbf{x}))$.

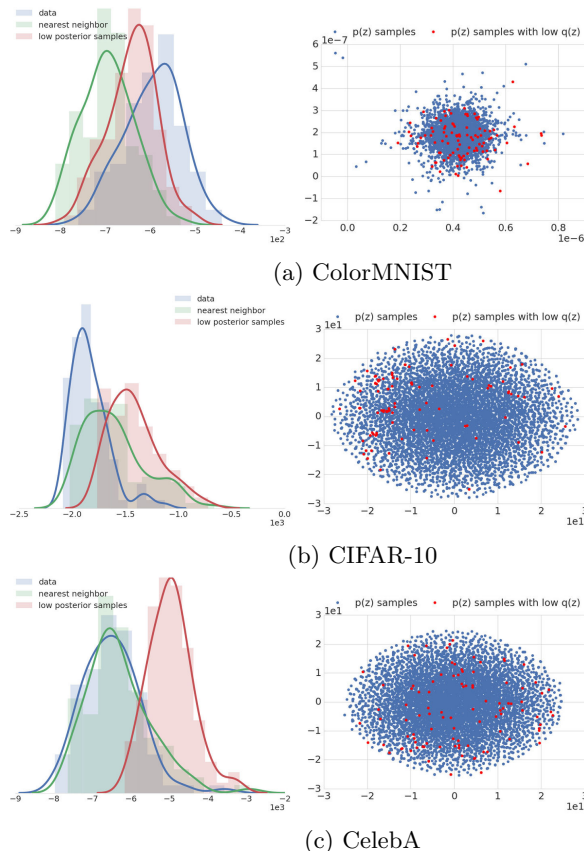


Figure 3: (Left) Evidence lower bound of uniformly sampled data, nearest neighbors to the low posterior samples from the dataset, and the low-probability posterior samples. We also plot the result of using a kernel density estimation for each histogram. (Right) 2D visualizations of the latent vectors used to create the VAE low posterior samples obtained using t-SNE.

Marginal distribution matching. Since conditional distribution matching fails to match marginal distributions in latent and visible space, directly minimizing $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ provides a compelling alternative. Adversarial Autoencoders (AAE) (Makhzani et al., 2015) ignore the mutual information term in Equation 5 and minimize $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ using adversarial training. AAEs are connected to marginal distribution matching not only by the ELBO, but also via Wasserstein distance (Tolstikhin et al., 2018) (see Appendix C).

4.2 Effects of density ratios on evaluation, scalability and latent representations

Evaluation. Using the density ratio trick in variational inference comes at a price: the loss of an estimate for the variational lower bound. Leveraging

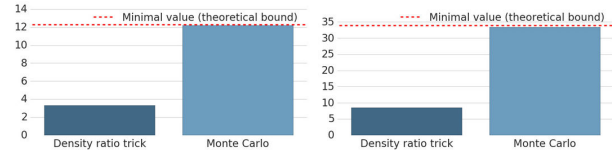


Figure 4: The density ratio estimator underestimates $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ for (left) Color-MNIST; (right) CIFAR.

binary classifiers to estimate KL divergences results in underestimated KL values, even when the discriminator is trained to optimality (see Figure 4). This issue is exacerbated in practice, as the discriminator is updated online during training. By using density ratios, VAE-GAN hybrids lose a meaningful estimate to report and a quantity to use to assess convergence, as shown in Figure 5. While VAEs can be compared with other likelihood models, VAE-GAN hybrids can only use data quality metrics for evaluation and lack a metric that can be used to track model convergence.

Scalability is another concern when using the density ratio trick. Using synthetic experiments where the true KL divergence is known, we show that using classifiers to provide learning signal does not scale with data dimensions (Appendix D). In practice this can be observed by training VAEs and AAEs with 1000 and 10000 latents. While VAEs learn to ignore extra latents (Figure 6) and scale with no architectural changes to extra dimensions, AAEs struggle to model ColorMNIST and completely fail on CelebA even with a bigger discriminator – see Figures 16 and 17 in Appendix B.

Latent representations learned via marginal divergence minimization have different properties than representations learned with VAEs. Figure 6 contrasts the posterior Gaussian means learned using AAEs and VAEs with 1000 latents on ColorMNIST: VAEs learn sparse representations, by using only a few latents to reconstruct the data and ignoring the rest, while AAEs learn dense representations. This is a consequence of marginal distribution matching: by not having a cost on the conditional $q_\eta(\mathbf{z}|\mathbf{x})$, AAEs can use any posterior latent distribution but lose the regularizing effect of $\text{KL}[q_\eta(\mathbf{z}|\mathbf{x}_n)||p(\mathbf{z})]$ - when learning a Gaussian $q_\eta(\mathbf{z}|\mathbf{x}_n)$ the variance collapses to 0 and the mean has a much wider range around the prior mean compared to VAEs. This finding suggests that for downstream applications such as semi supervised learning, using VAE representations might have a different effect compared to using representations learned by VAE-GAN hybrids.

5 GANS IN VISIBLE SPACE

The benefits and challenges of marginal distribution matching and implicit distributions in variational inference go beyond matching distributions in latent space.

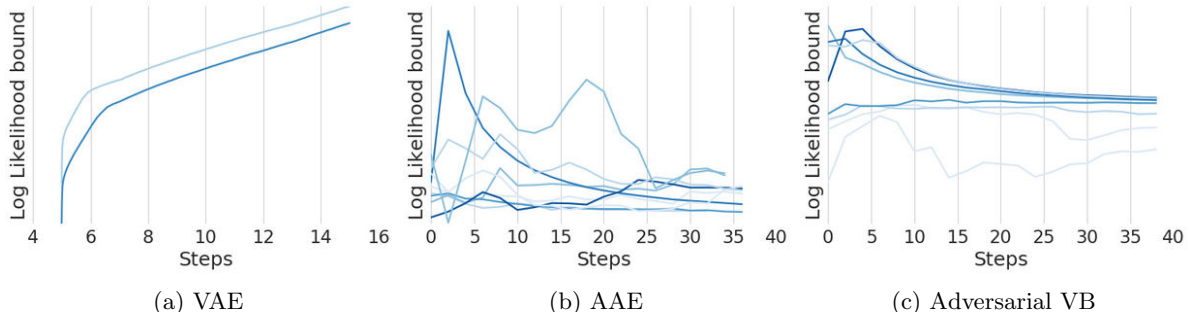


Figure 5: Variational lower bound estimates during training on ColorMNIST across hyperparameters. For accurate estimates, we expect that as model training progresses, the likelihood increases. Similar AdversarialVB results were showed in a maximum likelihood setting, see Figure 18 in Danihelka et al. (2017).

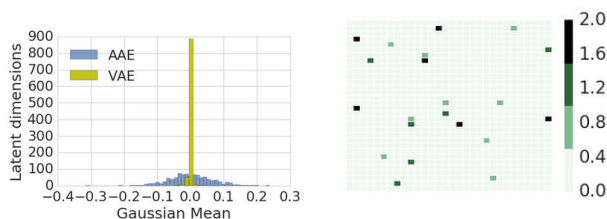


Figure 6: The effects of marginal distribution matching on latent representations. (Left) Learned latent means. (Right) Hinton diagram, showing VAE posterior KLs for dimension of the latent variable reshaped as a matrix for visualization, and showing that most KL values are low and close to zero.

Using *implicit models for the data* avoids notorious problems with explicit likelihoods, such as blurry samples and reconstructions produced by common choices such as Gaussian or Laplacian distributions. To introduce marginal distribution matching in variational inference, we introduce a ratio in the likelihood term of (4):

$$\begin{aligned} & \mathbb{E}_{p^*(\mathbf{x})} \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] \\ &= \mathbb{E}_{p^*(\mathbf{x})} \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})}{p^*(\mathbf{x})} \right] + \mathbb{H}[p^*(\mathbf{x})] \end{aligned} \quad (7)$$

$\mathbb{H}[p^*(\mathbf{x})]$ denotes the entropy of $p^*(\mathbf{x})$ and can be ignored for optimization purposes. This approach is a form of synthetic likelihood (Dutta et al., 2016, Wood, 2010).

Equation (7) shows the challenges of replacing the likelihood term with a density ratio:

1. Naively using the density ratio trick requires an infinite number of discriminators, one for each $p_\theta(\mathbf{x}|\mathbf{z})$ with $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}_n)$.
2. The model $p_\theta(\mathbf{x}|\mathbf{z})$ no longer receives gradients, since the parametric estimator replacing $\frac{p_\theta(\mathbf{x}|\mathbf{z})}{p^*(\mathbf{x})}$ is evaluated at data points, as given by the expectations in Equation 7. The model has been absorbed by the ratio estimator - by providing samples to it.

To avoid both problems and continue analyzing the effects of marginal distribution matching in variational inference, we leverage GANs to train the generative model. A previously proposed solution explores joint distribution matching (Donahue et al., 2017, Dumoulin et al., 2017, Li et al., 2017, Srivastava et al., 2017).

5.1 Marginal distribution matching

To bypass Equation (7), instead of using conditional distribution matching done when maximizing the ELBO, we use a ratio estimator to distinguish between $p^*(\mathbf{x})$ and the marginal reconstruction distribution $p_{\theta^R(\mathbf{x})}^R = \int p_\theta(\mathbf{x}|\mathbf{z})q_\eta(\mathbf{z})d\mathbf{z}$. The loss function for the ratio estimator $\mathcal{D}_\phi(\mathbf{x})$ becomes:

$$\mathbb{E}_{p^*(\mathbf{x})} [-\log \mathcal{D}_\phi(\mathbf{x})] + \mathbb{E}_{p_{\theta^R(\mathbf{x})}^R} [-\log(1 - \mathcal{D}_\phi(\mathbf{x}))] \quad (8)$$

To produce realistic reconstructions which fool the discriminator, the encoder and decoder minimize an adversarial loss, such as $-E_{p_{\theta^R(\mathbf{x})}^R} \log(1 - \mathcal{D}_\phi(\mathbf{x}))$. This approach to marginal distribution matching in visible space can be used in conjunction with either conditional or marginal matching in latent space and lends itself to training both implicit and explicit models.

When learning explicit models, marginal distribution matching and conditional distribution matching as done in variational inference can be combined. From the variational inference view the adversarial loss acts as a regularizer, steering reconstructions towards an area of the space that makes them realistic enough to fool the discriminator. From the adversarial perspective, the reconstruction loss can be viewed as a regularizer which avoids mode collapse, a notorious issue with GANs Arjovsky et al. (2017), Srivastava et al. (2017).

All VAE-GAN hybrids discussed so far rely on the KL term to match distributions in latent space. Adversarial training in visible space gives us another way to match $q_\eta(\mathbf{z})$ and $p(\mathbf{z})$, by matching the data distribu-

tion and the model distribution, via a discriminator which minimizes:

$$2 \mathbb{E}_{p^*(\mathbf{x})} [-\log \mathcal{D}_\phi(\mathbf{x})] + \mathbb{E}_{p_{\theta^R}(\mathbf{x})} [-\log(1 - \mathcal{D}_\phi(\mathbf{x}))] + \mathbb{E}_{p_\theta(\mathbf{x})} [-\log(1 - \mathcal{D}_\phi(\mathbf{x}))] \quad (9)$$

The decoder now receives an adversarial loss both for reconstructions and samples, learning how to produce compelling looking samples from early on in training, like GANs. Unlike GANs, this model is able to do inference, by learning $q_\eta(\mathbf{z}|\mathbf{x})$. We will later show that this loss helps improve sample quality compared to using an adversarial loss only on reconstructions, as described in equation 8.

5.2 Joint variants

Matching *joint* distributions matches both visible and latent space distributions. Minimizing divergences in the joint space leads to solutions for lack of gradients and challenges of density ratio estimation of Equation 7 while staying completely in the variational inference framework, by changing the goal of the model and using a different variational bound. VEEGAN Srivastava et al. (2017) changes the model objective from matching the marginals $p^*(\mathbf{x})$ and $p_\theta(\mathbf{x})$ in data space to matching marginal distributions in latent space. Using a reconstructor network to learn a posterior distribution $p_\gamma(\mathbf{z}|\mathbf{x})$, the goal of VEEGAN is to match $p_\gamma(\mathbf{z})$ to the prior $p(\mathbf{z})$. Because $p_\gamma(\mathbf{z})$ is intractable, $p_\theta(\mathbf{x}|\mathbf{z})$ is introduced as a variational distribution using the lower bound:

$$\int p_\gamma(\mathbf{z}|\mathbf{x})p^*(\mathbf{x})d\mathbf{x} \leq \text{KL}[p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})||p_\gamma(\mathbf{z}|\mathbf{x})p^*(\mathbf{x})] + C$$

where C is a constant. The model is trained to minimize the negative of the lower bound and a reconstruction loss in latent space. To estimate $\text{KL}[p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})||p_\gamma(\mathbf{z}|\mathbf{x})p^*(\mathbf{x})]$, VEEGAN uses the density ratio trick - and hence also loses an estimate for the bound - Appendix H. Since the expectation in the objective is not taken with respect to the data distribution $p^*(\mathbf{x})$ but the joint distribution $p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ decoder gradients follow from the variational objective.

6 EVALUATION

We compare marginal distribution matching using implicit distributions in variational inference provided by VAE-GAN hybrids with VAEs and GANs on image generation tasks. We perform an extensive comparison between (detailed in Figure 7):

- VAEs with diagonal Gaussian posteriors – the most widespread variational inference model.
- DCGAN and WGAN-GP to provide a baseline for VAE-GAN hybrids.

- AAE to exhibit the effect of replacing the analytical KL in VAEs with the density ratio trick.
- VEEGAN to evaluate the density ratio trick on joint distributions, and the effect of reconstruction losses in latent space rather than data space.
- VGH and VGH++, two Variational GAN Hybrids we introduce inspired by the analysis in Section 5.1. Like AAEs, these models use marginal distribution matching in latent space via the density ratio trick. In visible space, VGH uses a discriminator trained on reconstructions (Equation 8), and VGH++ uses a discriminator trained on both data and samples (Equation 9). Both use an l_1 reconstruction loss in data space. Details are provided in Appendix 5. Contrasting VGH and AAEs measures the effect of explicit distribution matching in visible space, while comparing VGH and VGH++ assesses the efficacy of the density ratio trick in latent space.

We do not compare with AdversarialVB, as it underperformed in our initial experiments (Appendix J). Other approaches to VAE-GAN hybrids include replacing reconstruction losses on pixels with reconstruction losses on discriminator features (Larsen et al., 2016), but they are outside the scope of this work.

We train models on ColorMNIST, CelebA and CIFAR-10 and complement visual inspection of samples with three metrics: Inception Score, sample diversity, independent Wasserstein critic. We do not use the ELBO as an evaluation metric, since we have shown VAE-GAN hybrids cannot estimate it reliably. Descriptions of the metrics, experimental details and samples are in Appendices I, K and M.

Figures 8 and 9 show that VAEs perform well on datasets that have less variability, such as ColorMNIST and CelebA, but are not able to capture the subtleties of CIFAR10. Marginal distribution matching in visible space improves generation quality, with VEEGAN, VGH++ and VGH performing better than VAEs and AAEs. However, VAE-GAN hybrids do not outperform GANs on image quality metrics and consistently exhibit a higher sensitivity to hyperparameters, caused by additional optimization components. VGH++ consistently outperforms VGH, showing that the density ratio trick has issues matching the marginal latent posterior to the prior and that matching marginal distributions in visible space explains the increased sample quality of VAE-GAN hybrids compared to VAEs.

By assessing sample quality and diversity, our experiments show that, at present, merging variational inference with implicit and marginal distribution matching does not provide a clear benefit.

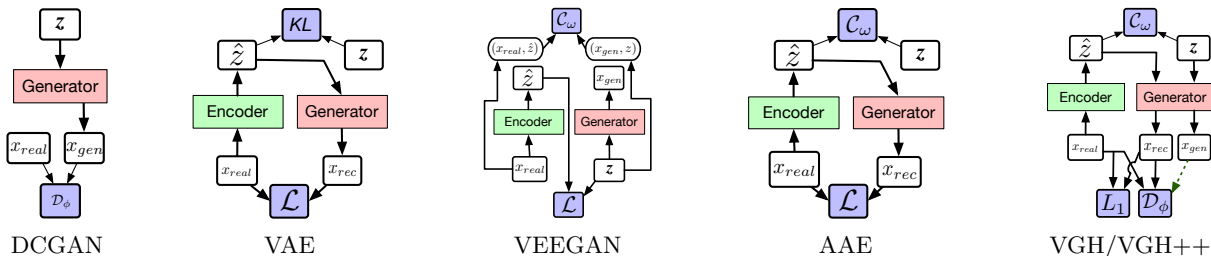


Figure 7: Model architectures (WGAN is similar to DCGAN). The difference between VGH++ and VGH is exemplified using the green arrow between x_{gen} and D_ϕ .

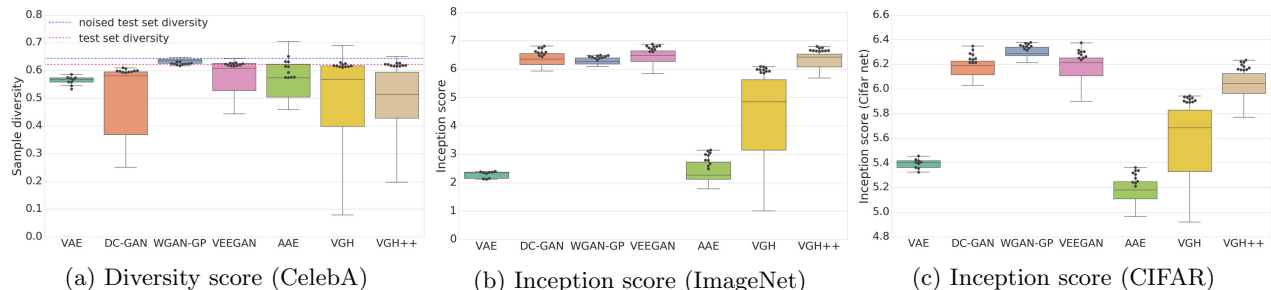


Figure 8: (Left) Sample diversity on CelebA, and is viewed relative to test set: too much diversity shows failure to capture the data distribution, too little is indicative of mode collapse. We also report the diversity obtained on a noised-version of the test set, which has a higher diversity than the test set. (Middle) Inception scores on CIFAR-10. (Right) Inception scores computed using a VGG-style network on CIFAR-10. For inception scores, higher values are better. For test data, diversity score: 0.621, inception score: 11.25, inception score (using CIFAR-10 trained net): 9.18. Best results are shown with black dots, and box plots show the hyperparameter sensitivity.

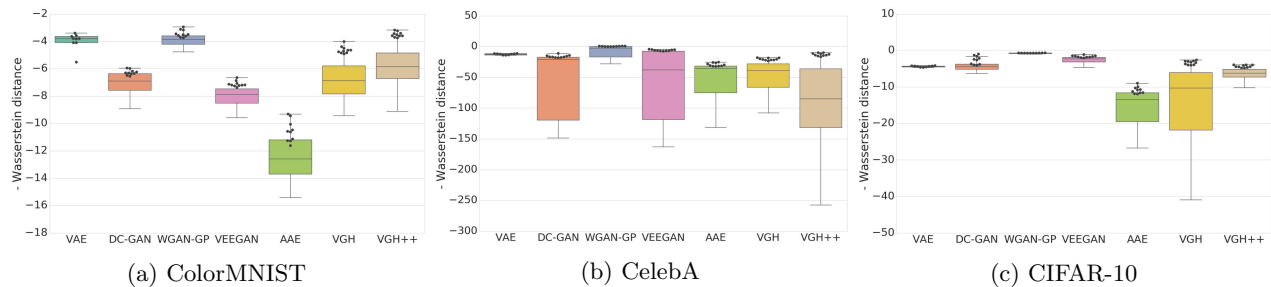


Figure 9: Comparison using negative Wasserstein distance computed using an independent Wasserstein critic. Higher is better. The metric captures overfitting and low quality samples and been shown to correlate with human evaluations (Jiwoong Im et al., 2018).

7 CONCLUSION

We have shown the widespread failure to learn marginal distributions with Variational Autoencoders. We asked whether this is the effect of conditional distribution matching, and of explicit posteriors and explicit model distributions. To test this hypothesis, we explored marginal distribution matching and implicit distribution in variational inference, through existing and new VAE-GAN hybrids.

Through a wide range of experiments, we have shown

that VAE-GAN hybrids do not deliver on the promise of addressing major challenges in variational inference. Problems with value estimation of divergences caused by the use of classifier probabilities, difficulties of scaling to high dimensions, hyperparameter sensitivity and the struggles to outperform GANs on sample quality metrics, limits, at present, the applicability of VAE-GAN hybrids. Since implicit models and adversarial training do not solve the obstacles of variational inference, distribution matching in latent and visible space remain important generative models research issues.

References

- A. A. Alemi, B. Poole, I. Fischer, J. V. Dillon, R. A. Saurous, and K. Murphy. An information-theoretic analysis of deep latent-variable models. *arXiv preprint arXiv:1711.00464*, 2017.
- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein GAN. In *ICML*, 2017.
- I. Danihelka, B. Lakshminarayanan, B. Uria, D. Wierstra, and P. Dayan. Comparison of Maximum Likelihood and GAN-based training of Real NVPs. *arXiv preprint arXiv:1705.05263*, 2017.
- L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using Real NVP. *arXiv preprint arXiv:1605.08803*, 2016.
- J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. In *ICLR*, 2017.
- V. Dumoulin, I. Belghazi, B. Poole, A. Lamb, M. Arjovsky, O. Mastropietro, and A. Courville. Adversarially learned inference. In *ICLR*, 2017.
- R. Dutta, J. Corander, S. Kaski, and M. U. Gutmann. Likelihood-free inference by penalised logistic regression. *arXiv preprint arXiv:1611.10242*, 2016.
- W. Fedus, M. Rosca, B. Lakshminarayanan, A. M. Dai, S. Mohamed, and I. Goodfellow. Many paths to equilibrium: GANs do not need to decrease a divergence at every step. In *ICLR*, 2018.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- I. J. Goodfellow. On distinguishability criteria for estimating generative models. *arXiv preprint arXiv:1412.6515*, 2014.
- I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville. Improved training of Wasserstein GANs. In *NIPS*, 2017.
- C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In *ICML*, 2017.
- I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. β -VAE: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017a.
- I. Higgins, A. Pal, A. A. Rusu, L. Matthey, C. P. Burgess, A. Pritzel, M. Botvinick, C. Blundell, and A. Lerchner. DARLA: Improving zero-shot transfer in reinforcement learning. In *ICML*, 2017b.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- M. D. Hoffman and M. J. Johnson. ELBO surgery: yet another way to carve up the variational evidence lower bound. In *Workshop in Advances in Approximate Bayesian Inference, NIPS*, 2016.
- S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- D. Jiwoong Im, A. He Ma, G. W. Taylor, and K. Branson. Quantitatively evaluating GANs with divergences proposed for training. *ICLR*, 2018.
- D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- D. P. Kingma and M. Welling. Auto-encoding variational Bayes. In *ICLR*, 2013.
- D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *NIPS*, 2014.
- J. Kos, I. Fischer, and D. Song. Adversarial examples for generative models. *arXiv preprint arXiv:1702.06832*, 2017.
- A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. In *ICML*, 2016.
- C. Li, H. Liu, C. Chen, Y. Pu, L. Chen, R. Henao, and L. Carin. Alice: Towards understanding adversarial learning for joint distribution matching. In *Advances in Neural Information Processing Systems*, pages 5501–5509, 2017.
- Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *ICCV*, 2015.
- A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML*, 2013.
- A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.
- M. F. Mathieu, J. J. Zhao, J. Zhao, A. Ramesh, P. Sprechmann, and Y. LeCun. Disentangling factors of variation in deep representation using adversarial training. In *NIPS*, 2016.
- L. Mescheder, S. Nowozin, and A. Geiger. Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks. In *ICML*, 2017.
- L. Metz, B. Poole, D. Pfau, and J. Sohl-Dickstein. Unrolled generative adversarial networks. In *ICLR*, 2017.
- S. Mohamed and B. Lakshminarayanan. Learning in implicit generative models. *arXiv preprint arXiv:1610.03483*, 2016.
- A. Nguyen, J. Yosinski, Y. Bengio, A. Dosovitskiy, and J. Clune. Plug & play generative networks: Conditional iterative generation of images in latent space. *arXiv preprint arXiv:1612.00005*, 2016.
- A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier GANs. *arXiv preprint arXiv:1610.09585*, 2016.
- G. Papamakarios, I. Murray, and T. Pavlakou. Masked autoregressive flow for density estimation. In *NIPS*, 2017.

- Y. Pu, W. Wang, R. Henao, L. Chen, Z. Gan, C. Li, and L. Carin. Adversarial symmetric variational autoencoder. In *NIPS*, 2017.
- D. Rezende, I. Danihelka, K. Gregor, D. Wierstra, et al. One-shot generalization in deep generative models. In *ICML*, 2016.
- D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014.
- M. Rosca, B. Lakshminarayanan, D. Warde-Farley, and S. Mohamed. Variational approaches for auto-encoding generative adversarial networks. *arXiv preprint arXiv:1706.04987*, 2017.
- T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training GANs. In *NIPS*, 2017.
- L. I. Smith. A tutorial on principal components analysis. Technical report, 2002.
- A. Srivastava, L. Valkov, C. Russell, M. Gutmann, and C. Sutton. VEEGAN: Reducing mode collapse in GANs using implicit variational learning. In *NIPS*, 2017.
- M. Sugiyama, T. Suzuki, and T. Kanamori. Density-ratio matching under the Bregman divergence: a unified framework of density-ratio estimation. *Annals of the Institute of Statistical Mathematics*, 64(5):1009–1044, 2012.
- C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.
- L. Theis, A. van den Oord, and M. Bethge. A note on the evaluation of generative models. In *ICLR*, 2016.
- I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf. Wasserstein auto-encoders. In *ICLR*, 2018.
- Z. Wang, E. P. Simoncelli, and A. C. Bovik. Multiscale structural similarity for image quality assessment. In *Conference Record of the Thirty-Seventh Asilomar Conference on Signals, Systems and Computers, 2004.*, volume 2, pages 1398–1402. IEEE, 2003.
- S. N. Wood. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310):1102–1104, 2010.

Appendix – Distribution Matching in Variational Inference

A Low posterior VAE samples

In this section we detail the experiments performed to create low posterior VAE samples. For the algorithm used to obtain the low posterior samples, see Algorithm 1. For the low posterior samples alongside standard samples from the model, see Figure 10.

Algorithm 1 Pseudocode for generating low posterior VAE samples

```

1: Load trained variational model with prior  $p(\mathbf{z})$ , and
   posterior  $q_\eta(\mathbf{z}|\mathbf{x})$ .
2:  $\text{log\_q\_to\_z} = \{\}$ 
3: for  $i = 1 : \text{num\_z}$  do
4:   sample  $\mathbf{z}_i$  from  $p(\mathbf{z})$ 
5:    $\text{posterior\_list} = []$ 
6:   for  $\mathbf{x}$  in dataset do
7:     append  $\log q_\eta(\mathbf{z}|\mathbf{x})(\mathbf{z}_i)$  to  $\text{posterior\_list}$ 
8:   end for
9:    $\log q(\mathbf{z}_i) = \text{log\_mean\_exp}(\text{posterior\_list})$ 
10:   $\text{log\_q\_to\_z}[\log q(\mathbf{z}_i)] = \mathbf{z}_i$ 
11: end for
12: Sort  $\text{log\_q\_to\_z}$  by key and let  $\mathbf{z\_adv}$  be the list of
   values corresponding to the  $n$  smallest keys.
13:  $\text{low\_posterior\_samples} = []$ 
14: for  $\mathbf{z}$  in  $\mathbf{z\_adv}$  do
15:   append a sample from  $p_\theta(\mathbf{x}|\mathbf{z})$  to
    $\text{low\_posterior\_samples}$ 
16: end for
17: return  $\text{low\_posterior\_samples}$ 

```

A.1 Low posterior samples - model analysis

In order to understand why the low posterior samples look as shown in Figure 10, we performed an analysis to show how these samples compare to the data distribution. For ColorMNIST, we visually saw that the samples are thicker than the data or the standard samples, while for CelebA and CIFAR-10 we saw predominantly white backgrounds so we plotted the pixel histogram of the dataset against the pixel histogram of the low posterior samples (Figure 12). The histograms of pixels for data, low posterior samples and their nearest neighbors in the dataset shows that the low posterior samples differ in pixel composition compared to the uniformly sampled data, but to check whether images like this exist in the dataset we plot the nearest neighbors in l_2 distance from the dataset to the low posterior samples (see Figure 11). This analysis shows that data examples that are similar to the low posterior samples exist, but based on the histogram analysis and visual inspection

we know that they have low probability under the true data distribution. A hypothesis emerges: the model is not putting mass in the marginal posterior distribution for areas of the space that encode data points which have low probability under the true data distribution. To test this hypothesis, we report the histograms of average $\text{KL}[q_\eta(\mathbf{z}|\mathbf{x}_n)||p(\mathbf{z})]$, obtained by encoding data sampled from the dataset uniformly compared to the nearest neighbors of the low posterior samples Engineering Retreat Google examples and the low posterior samples themselves (Figure 13). Our results show that indeed, the data points closest to the low posterior samples have a higher KL cost compared to datapoints sampled uniformly from the dataset and that these data points are unlikely under the data distribution. Knowing that these examples are unlikely under the true data distribution, we expect to see the same under the *model* distribution. In Figure 14 we show that for CIFAR-10 and CelebA, the model reports that the low posterior samples are more likely than the data. This demonstrated that the model is unable to capture the subtleties of the data distribution, and can be fooled into predicting high likelihoods for samples that have low probability under the true data distribution. In this work we have exploited the gap between the prior and marginal posterior distributions in VAEs trained with Gaussian posteriors, to show that the model can generate samples far from the sample distribution and far from the data distribution. Previous work has focused on finding adversarial examples as input to the VAE by finding points in data space that the VAE is unable to reconstruct (Kos et al., 2017).

B Scaling up latent spaces in VAEs and AAEs

Figure 16 shows that VAEs scale better than AAEs to a high number of latent size - in this case 10000.

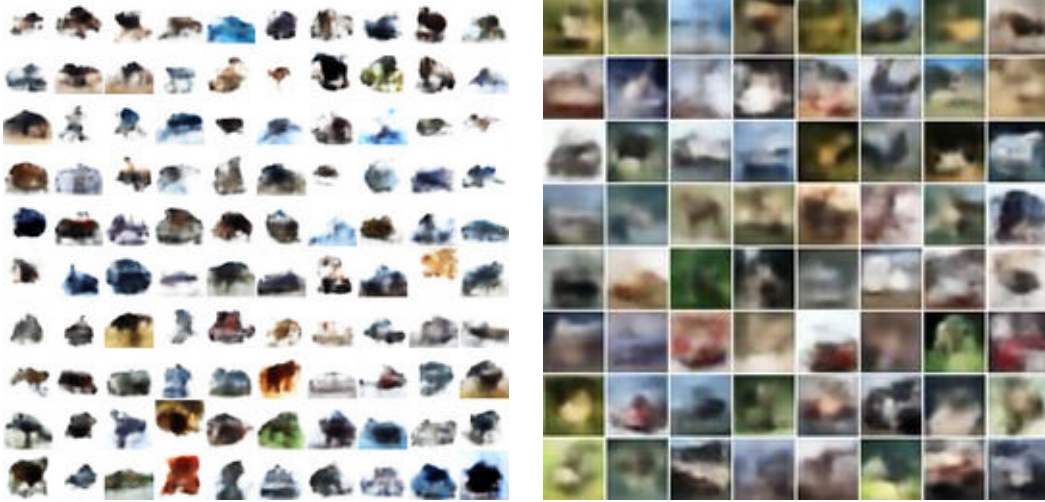
C Connection to Wasserstein Autoencoders

Tolstikhin et al. (2018) prove the connection between using optimal transport to minimize the distance the true data distribution and the model distribution induced by an *implicit* model and using autoencoders which minimize a distance between the marginal posterior and prior in latent space. Specifically, they show that:

Theorem 1 *Let c be a measurable cost function with values in R_+ and $\mathbb{P}(p^*(\mathbf{x}), p_\theta(\mathbf{x}))$ the set of all joint distributions with marginals $p^*(\mathbf{x})$ and $p_\theta(\mathbf{x})$, respectively. For a model where $p_\theta(\mathbf{x}|\mathbf{z})$ is a Dirac delta function,*



(a) ColorMNIST

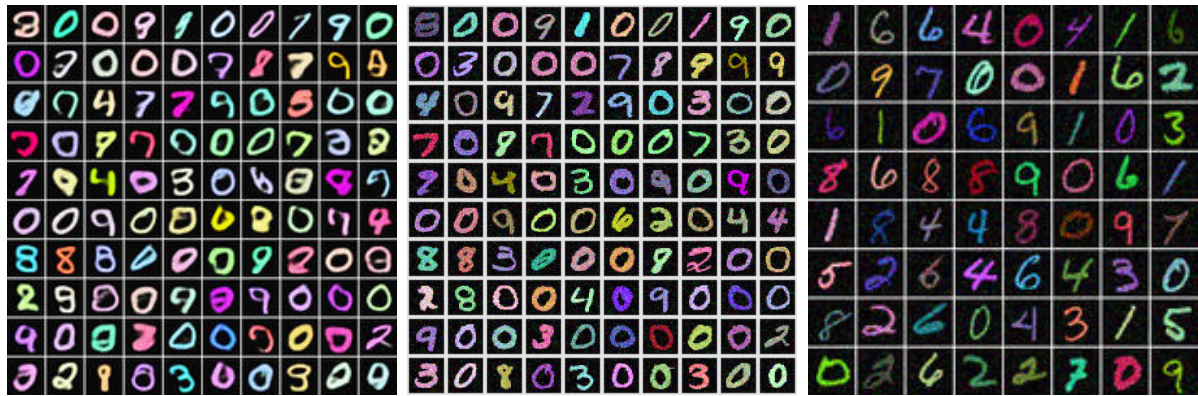


(b) CIFAR-10

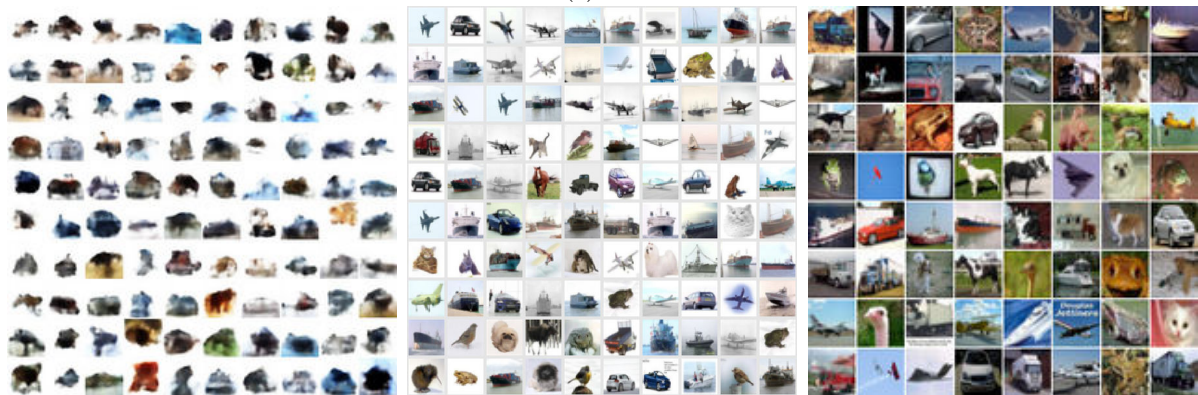


(c) CelebA

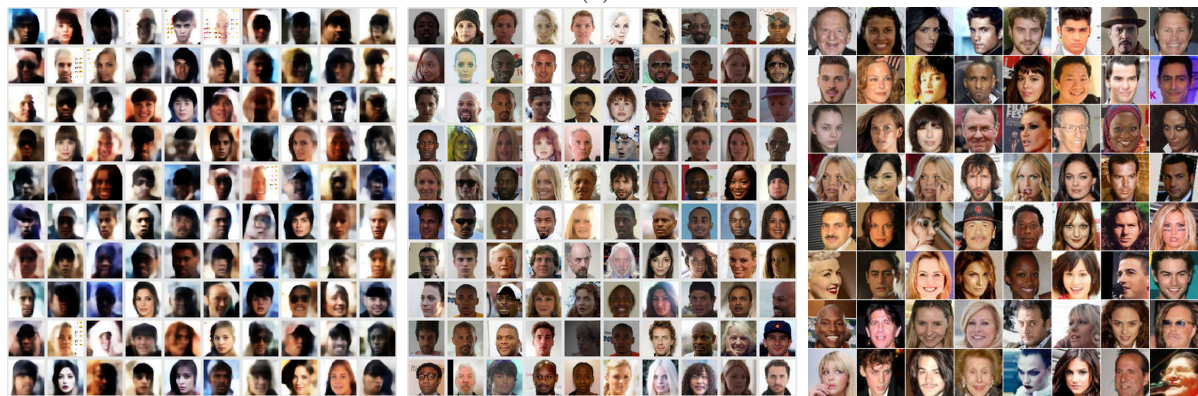
Figure 10: Low posterior VAE samples obtained by sampling latent variables from an area in the latent space where the marginal posterior has low probability mass (left), alongside standard VAE samples from the **same** model (right).



(a) ColorMNIST



(b) CIFAR-10



(c) CelebA

Figure 11: Low posterior samples (left), the nearest neighbors in the dataset from the low posterior samples (middle), uniformly dataset examples (right).

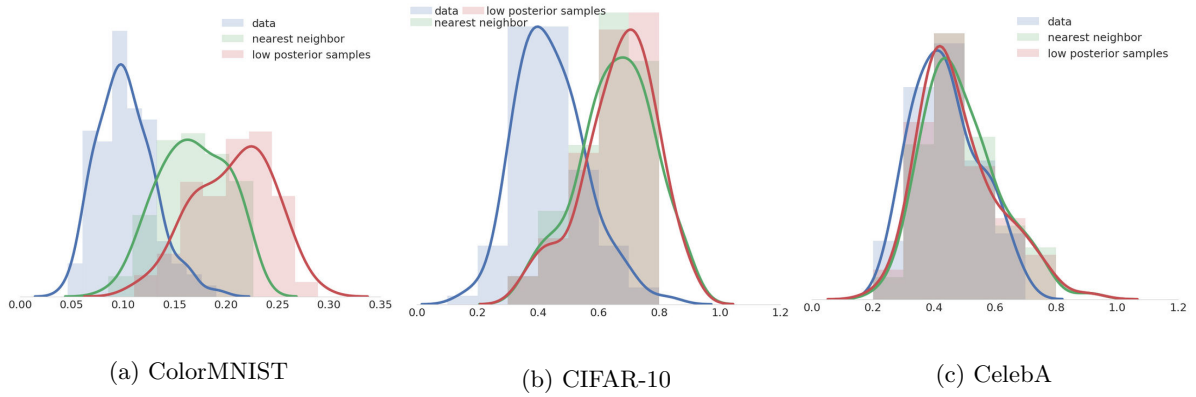


Figure 12: Histogram of pixels on uniformly sampled data, the nearest neighbors from the low posterior samples and the low posterior samples. For CIFAR-10 and ColorMNIST we see that both the low posterior samples and their nearest neighbors in the dataset are atypical compared to the data. We also plot the result of using a KDE density estimation for each histogram.

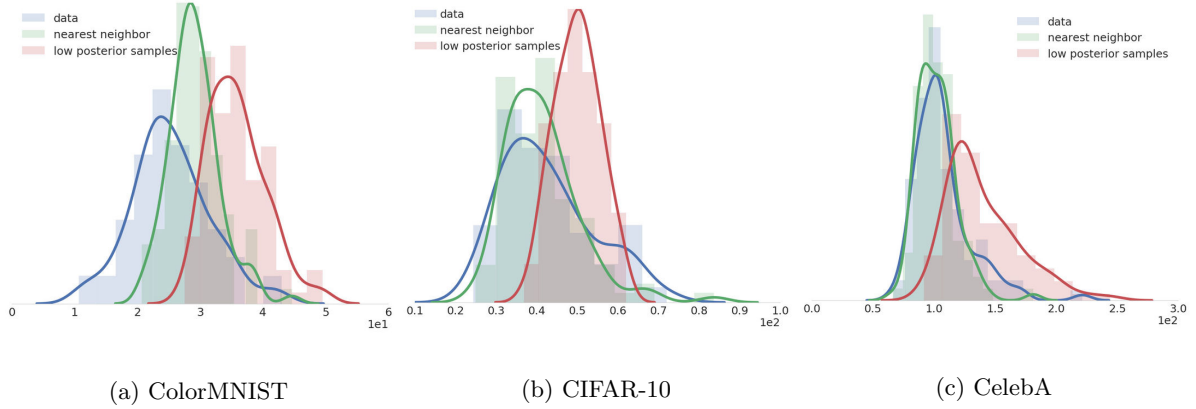


Figure 13: Histogram of $KL[q_\eta(\mathbf{z}|\mathbf{x}_n)||p(\mathbf{z})]$ on uniformly sampled data, the nearest neighbors from the low posterior samples and the low posterior samples. Overall, we see a higher KL term for data points close to the low posterior samples and for low posterior samples. We also plot the result of using a KDE density estimation for each histogram.

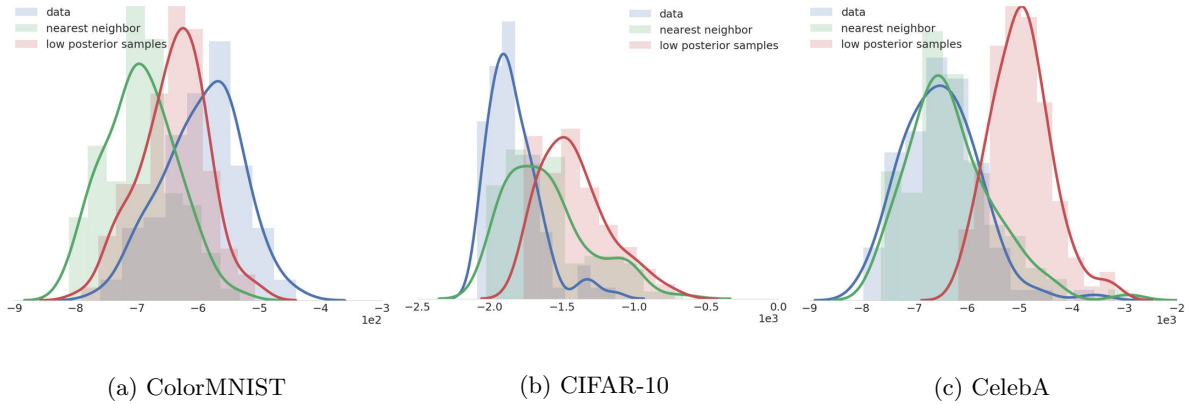


Figure 14: Model evidence lower bound of uniformly sampled data, the nearest neighbors from the low posterior samples in the dataset and the low posterior samples. While on ColorMNIST the model recognizes that the low posterior samples and their nearest neighbors have a lower likelihood than the data, for CIFAR-10 and CelebA the model thinks the low posterior samples are more likely than the data. We also plot the result of using a KDE density estimation for each histogram.

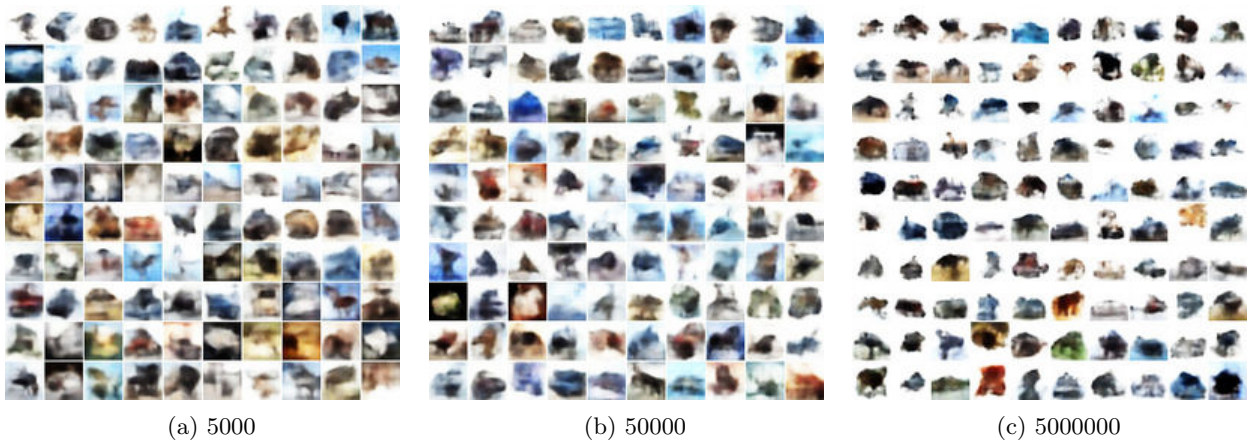


Figure 15: Low posterior samples on CIFAR-10 resulting from different numbers of latents sampled from the prior (corresponding to num_z in Algorithm 1). While the samples get more pathological with an increased number of samples from the prior, we can already generate abnormal VAE samples from a small number of latent samples.

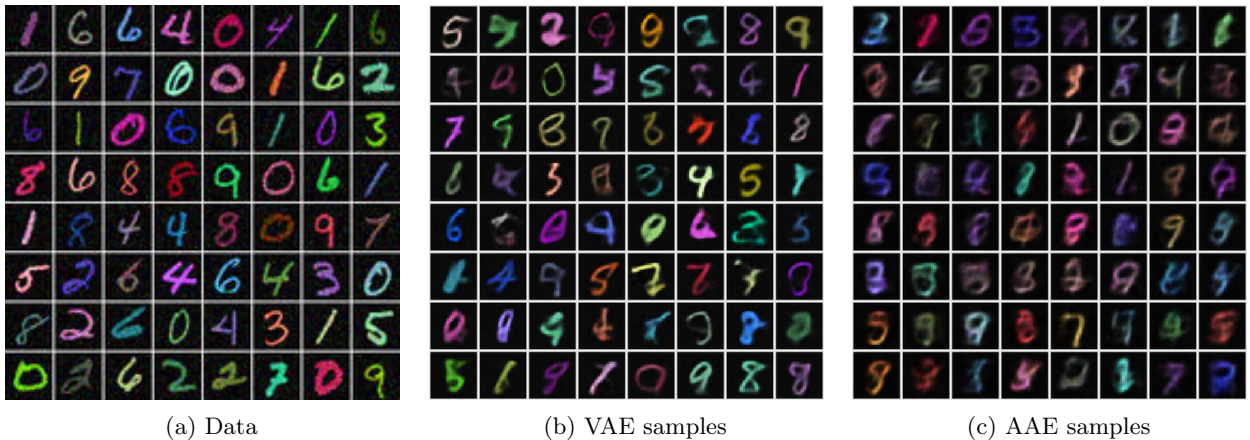


Figure 16: Samples from VAE and AAE with 10k latents compared to data on ColorMNIST.

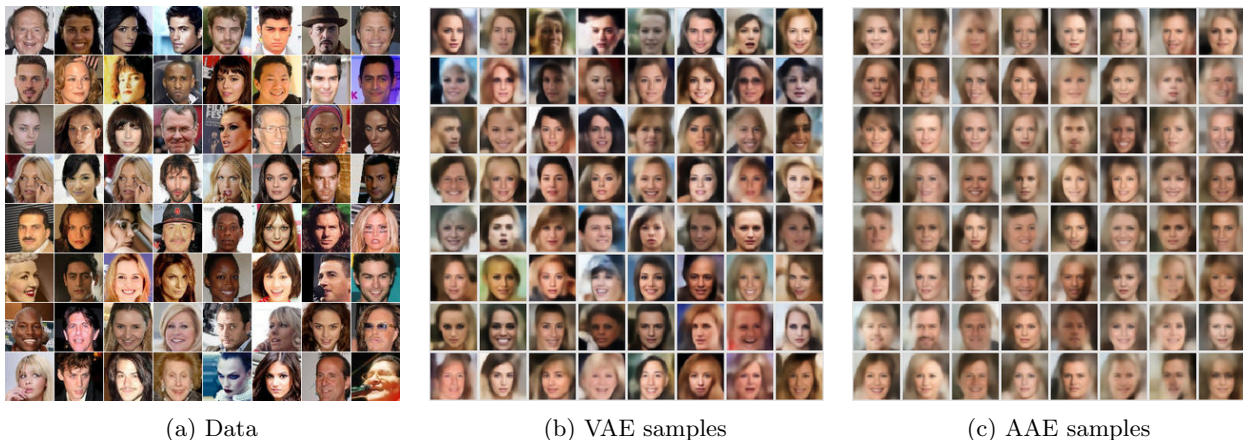


Figure 17: Samples from VAE and AAE with 8k latents compared to data on CelebA.

ie. \mathbf{z} is mapped to \mathbf{x} deterministically $\mathbf{x} = G(\mathbf{z})$, the following holds:

$$\begin{aligned} & \inf_{\Gamma \sim \mathbb{P}(p^*(\mathbf{x}), p_{\theta}(\mathbf{x}))} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \Gamma} [c(\mathbf{x}, \mathbf{y})] \\ &= \inf_{q_{\eta}(\mathbf{z}|\mathbf{x}): q_{\eta}(\mathbf{z})=p(\mathbf{z})} \mathbb{E}_{\mathbf{x} \sim p^*(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_{\eta}(\mathbf{z}|\mathbf{x})} [c(\mathbf{x}, G(\mathbf{z}))] \end{aligned} \quad (10)$$

Theorem 1 shows that for implicit models, encoder-decoder models can be introduced as a way to make the optimal transport computation tractable, as optimizing over the space of joint distributions $\mathbb{P}(p^*(\mathbf{x}), p_{\theta}(\mathbf{x}))$ is not feasible - another approach to make the computation tractable, when c is a metric, is to use the Kantorovich-Rubinstein duality Arjovsky et al. (2017).

Similarly, when doing maximum likelihood - minimizing the KL divergence between the model and data distribution - encoder-decoder models are introduced to overcome the intractability introduced by Equation 1, using Jensen’s inequality (Equation 3).

The connections between optimal transport and maximum likelihood do not stop here - Tolstikhin et al. (2018) optimize the RHS of Equation 10 using a relaxation that adds a penalty to the objective which forces the marginal $q_{\eta}(\mathbf{z})$ close to $p(\mathbf{z})$:

$$\mathbb{E}_{\mathbf{x} \sim p^*(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q_{\eta}(\mathbf{z}|\mathbf{x})} [c(\mathbf{x}, G(\mathbf{z}))] + \lambda D_{\mathbf{z}}(q_{\eta}(\mathbf{z}), p(\mathbf{z})) \quad (11)$$

When setting c to the l_2 or l_1 distance (corresponding to a Gaussian or Laplacian likelihood in the explicit model case), and setting $D_{\mathbf{z}}$ to the KL divergence, we obtain the bound obtained by doing the ELBO surgery performed on the maximum likelihood objective and the training objective used to train Adversarial Autoencoders.

Hence under certain conditions, optimal transport and maximum likelihood problems lead to encoder-decoder architecture and similar optimization criteria, but allow for different modeling choices. In variational inference, the choice of $p_{\theta}(\mathbf{x}|\mathbf{z})$ in decides the reconstruction cost function, while in Wasserstein Autoencoders $p_{\theta}(\mathbf{x}|\mathbf{z})$ has to be a Dirac delta function, but the reconstruction cost is chosen by the practitioner. Like Adversarial Autoencoders, by approximating $D_{\mathbf{z}}$ using the density ratio trick or the Wasserstein GAN objective, Wasserstein Autoencoders lose a meaningful quantity to track - which is not the case for Variational Autoencoders which use the closed form of the KL divergence.

The connection between optimal transport and maximum likelihood opens new avenues for research that we leave for future work, while the different modeling choices provide new flexibility to machine learning practitioners.

D Learning with density ratios: synthetic data experiments

We show that the density ratio trick can be used for learning, even though it cannot be used for estimating divergences. We also show the challenges of scaling the density ratio trick to higher dimensions. To do so, devise a set of synthetic experiments where the true KL divergence is known and where we can determine how this approach scales with data dimensionality.

We use Gaussian distributions, defined by passing a random normal vector through an affine transformation with $\mathbf{W} \in \mathbb{R}^{kd}$, $\mathbf{b} \in \mathbb{R}^d$, $\mathbf{z} \in \mathbb{R}^k$:

$$\mathbf{z} \sim \mathcal{N}(0, \mathbb{I}_k), \mathbf{x} = \mathbf{W}^{\top} \mathbf{z} + \mathbf{b} \implies x \sim \mathcal{N}(\mathbf{b}, \mathbf{W}^{\top} \mathbf{W}) \quad (12)$$

To ensure $\mathbf{W}^{\top} \mathbf{W}$ is full rank, we set $d = k/10$ in all our experiments.

We first train a classifier to distinguish between two such Gaussian distributions, for varying values of d . In the first setting, we are concerned with *divergence estimation*, keeping both distributions fixed and only learn the density ratio using the classifier. Once the classifier is trained, we report the difference between the estimated and true KL divergence values. In the second setting, we are concerned with *divergence minimization* and we begin with the same initialization for the two distributions, but learn the parameters of the second Gaussian (\mathbf{W} and \mathbf{b} in Equation 12) to minimize the estimated divergence between the two distributions. This is a GAN training regime, where the generator loss is given by the reverse KL generator loss (Mohamed and Lakshminarayanan, 2016): $-\log \frac{D(\mathbf{x})}{1-D(\mathbf{x})}$. We track the true KL divergence during training together with the online classifier estimated divergence - we should not expect the latter to be accurate, the classifier is not trained to optimality for each update of the learned distribution, as the two models are trained jointly.

If for the same classifier that failed to approximate the true KL divergence in the estimation experiments we observe a decrease in true divergence in the learning experiments, we can conclude that while the density ratio trick might not be a useful for estimation, it can still be used as an optimization tool. To ensure our conclusions are valid, we control over hyper-parameters, classifier architectures and random seeds of the Gaussian distributions, and average results over 10 runs.

Our main findings are summarized in Figures 18, 19, and 20 and reveal that using density ratios for learning does not reliably scale with data dimensionality. For lower dimensional data (1 and 10 dimensions), the model is able to decrease the true KL divergence. However, for higher dimensional Gaussians (dimensions 100 and 1000), a classifier with 100 million parameters (4 layer MLP) is not able to provide useful gradients and learning diverges (rightmost plot in Figure 2). Regardless of data dimensionality, the estimate of the true KL divergence provided by the density ratio trick was not accurate.

The discriminator was trained with the AdamOptimizer with β_1 set to 0.5 and β_2 set to 0.9 for 1000000 iterations. Unless otherwise specified, the discriminator was a 4 layer MLP, trained with a learning rate of 0.0001. The learning rate used for learning the Gaussian was 0.001. Similar results were obtained for different learning rates for the discriminator and the linear model of the Gaussian distribution.

E Estimating $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$

We present the details of marginal KL estimation experiments described in Sections 3 and 4.

E.1 Estimating $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ using the Monte Carlo approach

Algorithm 2 describes the approach used to estimate $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ via Monte Carlo methods. While this approach is computationally expensive, it the most accurate one. We now describe the details of this computation. For each \mathbf{x}_i we used 10^6 samples from $q(\mathbf{z}|\mathbf{x})$ to estimate $\log \frac{q_\eta(\mathbf{z})}{p(\mathbf{z})}$. To estimate $q_\eta(\mathbf{z})$ for latent sample we used the entire dataset training and validation split of the dataset at hand. In all our figures and tables, this number is reported as N .

Algorithm 2 Pseudocode for estimating the marginal KL using MC

- 1: Load trained variational model with prior $p(\mathbf{z})$, and posterior $q_\eta(\mathbf{z}|\mathbf{x})$.
 - 2: `marginal_kl = 0.0`
 - 3: **for** `i = 1 : num_z` **do**
 - 4: sample \mathbf{x}_i from $p^*(\mathbf{x})$, sample z_i from $q_\eta(\mathbf{z}|\mathbf{x}_i)$
 - 5: `posterior_list = []`
 - 6: **for** `x` in dataset **do**
 - 7: append $\log q_\eta(\mathbf{z}|\mathbf{x})(z_i)$ to `posterior_list`
 - 8: **end for**
 - 9: `log q(zi) = log_mean_exp(posterior_list)`
 - 10: `marginal_kl += log q(zi) - log p(zi)`
 - 11: **end for**
 - 12: `marginal_kl = marginal_kl/num_z`
-

E.2 Estimating $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ using the density ratio trick

To estimate $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ using the density ratio trick as shown in Figure 4 we used Algorithm 3. For all datasets, we noticed that this approach is highly sensitive to hyperparameters. We explain this two fold: first, this approach relies on the probabilities reported by a neural network classifier, which have been known to be inaccurate. New methods have been proposed to address this issue (Guo et al., 2017), and we leave exploring these approaches for future work. Second, as shown in Section E.3, the distribution $q(z)$ can be very complex, making it hard for the classifier to learn to distinguish between samples from the two distributions.

We show the hyperparameter sensitivity by training different models to estimate the marginal $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ for the same VAE and report the different values obtained. All trained density ratio estimators were MLPs with Leaky Rectified activations of slope 0.2 and were trained for $5 * 10 * 5$ steps using the AdamOptimizer with β_1 and β_2 equal to 0.9.

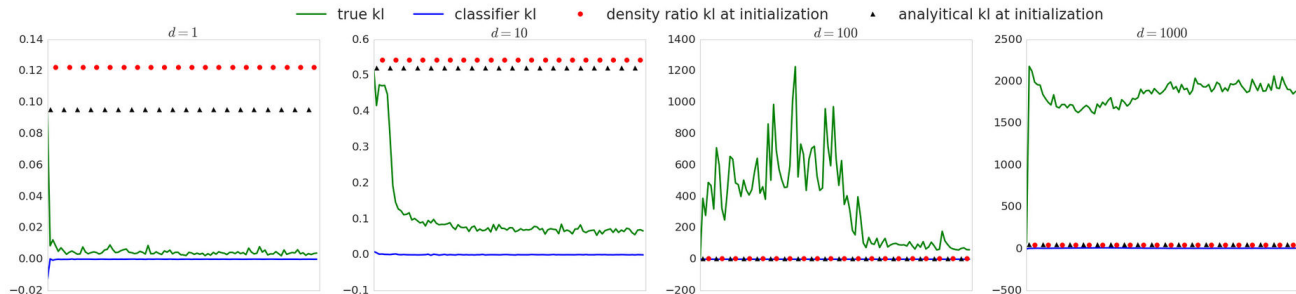


Figure 18: Divergence estimation and minimization of Gaussian distributions, for different data dimensions d . We plot training progress using the true KL divergence between the learned and true distributions. As a reference point, we show the true KL divergence at initialization, together with how well the same classifier architecture is able to estimate the initial true KL when the two Gaussian distributions are stationary. Results are averaged over 10 different initializations for the classifier.



Figure 19: Synthetic Gaussian experiments, for data different dimensions d , with a higher learning rate for the discriminator: 0.001. The model is more unstable, no longer being able to converge for data of dimensionality 100. Results are averaged over 10 different initializations for the discriminator.

Algorithm 3 Pseudocode for estimating the marginal KL using the density ratio trick

- 1: Load trained variational model with prior $p(\mathbf{z})$, and posterior $q_\eta(\mathbf{z}|\mathbf{x})$.
 - 2: Initialize code discriminator parameters ω randomly.
 - 3: **for** iter = 1 : max_iter **do**
 - 4: Update parameters ω by maximizing $\mathbb{E}_{p^*(\mathbf{x})}\mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})}[\log(\mathcal{C}_\omega(\mathbf{z}))] + \mathbb{E}_{p(\mathbf{z})}[\log(1 - \mathcal{C}_\omega(\mathbf{z}))]$
 - 5: **end for**
 - 6: marginal_kl = 0.0
 - 7: **for** i = 1 : num_z **do**
 - 8: sample \mathbf{x}_i from $p^*(\mathbf{x})$, sample \mathbf{z}_i from $q_\eta(\mathbf{z}|\mathbf{x}_i)$
 - 9: marginal_kl += $\log \mathcal{C}_\omega(\mathbf{z}_i) - \log(1 - \mathcal{C}_\omega(\mathbf{z}_i))$
 - 10: **end for**
 - 11: marginal_kl = marginal_kl/num_z
-

Table 1: Estimating a marginal KL using the density ratio trick for a standard VAE with 50 latents trained on Color MNIST. The number of hidden units per layer was 5000. When the KL is estimated numerically, the result is 12.3. From Equation (5) and that the mutual information term is bound by $\log N$, with $N = 60000$ and the average posterior KL of the model is 23.34, we know that the value needs to be greater than 12.46. All models used a learning rate of 0.0005.

# layers	Gradient Penalty	Activation Noise	KL
3	No	No	2.3
4	No	No	3.3
4	No	Yes	25812.8
4	Yes	No	3.1

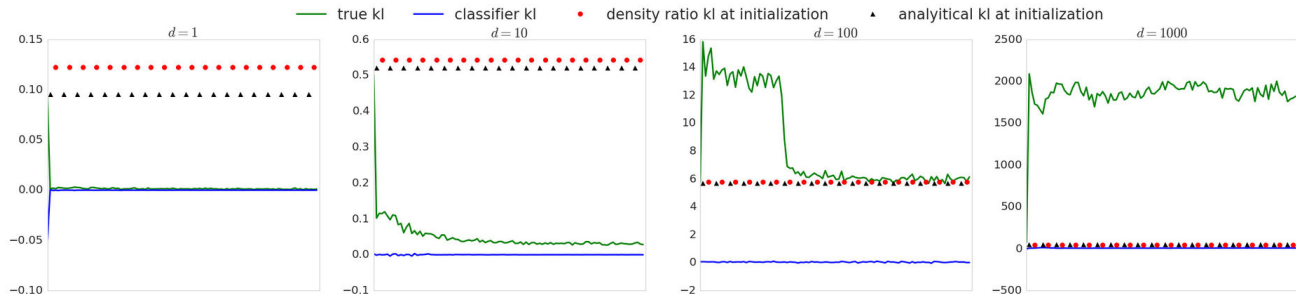


Figure 20: Synthetic Gaussian experiments, for data different dimensions d , with 5 discriminator updates. The experimental set up is exactly the same as Figure 18, including random seeds and hyperparameters but the discriminator is updated 5 times per generator update. Making the discriminator updates more frequent makes the learned model converge earlier for $d = 1, 10, 100$ but no improvement for $d = 1000$. We see no improvement in the estimated KL, even in cases where the discriminator could estimate the KL when trained to optimality, reported as 'density ratio KL as initialization'.

Table 2: Estimating a marginal KL using the density ratio trick for a standard VAE with 100 latents trained on CelebA. The number of hidden units per layer was 5000. When the KL is estimated numerically, the result is 100.3. From Equation (5) and that the mutual information term is bound by $\log N$, with $N = 162770$ and the average posterior KL of the model is 112.37, we know that the value needs to be greater than 100.0.

# layers	Learning Rate	Activation Noise	Dropout	KL
5	0.0005	No	No	17.7
5	0.0005	Yes	No	720140.9
5	0.0001	No	No	14.8
7	0.0005	No	No	18.8
7	0.0005	No	Yes	19.0
7	0.0001	No	No	18.0
10	0.0001	No	No	18.1

Results on ColorMNIST are summarized in Table 1, while CelebA results are summarized in Table 2.

While analyzing these results, we observed that adding noise to the activation of the classifier resulted in a better classifier, but also a more confident one, which underestimates the probability that a sample was given by the prior, and the KL value being over estimated. We also see that the resulting KL is quite sensitive to the the architecture of the classifier, with an extra layer resulting in a substantial value increase for the Color MNIST case. Gradients penalties and dropout did not result in a big change in the estimated value.

E.3 Estimating $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$ using a density model for $q(\mathbf{z})$

We compare the MonteCarlo and density ratio trick approach with a third way to estimate $\text{KL}[q_\eta(\mathbf{z})||p(\mathbf{z})]$:

by using a density model to learn $q(\mathbf{z})$. We use three density models: Gaussian Mixture Model, Masked Auto-regressive Flows (Papamakarios et al., 2017), and Gaussian auto-regressive models implemented using an (Hochreiter and Schmidhuber, 1997). Algorithm 4 was used to learn these models. Across models, we show a failure (see Figure 21) to reach the theoretical bound, showing that $q_\eta(\mathbf{z})$ is a complex distribution.

Diagonal Gaussian Mixture Models

In this setting, we model $q_\eta(\mathbf{z})$ using:

$$q_\eta(\mathbf{z}) = \sum_i^k \pi_i \mathcal{N}(\mathbf{z}|\mu_i, \sigma_i), \sum_i \pi_i = 1$$

Masked Auto-regressive Flows

In this setting, the density model is given by transforming a standard Gaussian using auto-regressive models as normalizing flows:

$$q_\eta(\mathbf{z}) = \mathcal{N}(0, \mathbb{I}|\mathbf{z}) \left| \det \left(\frac{d f^{-1}}{dz} \right) \right|$$

where f has to be an invertible function for which the determinant of its Jacobian is easy to compute. In practice, we leverage the fact that the composition of two functions which have these proprieties also has this property to chain a number of transforms.

Gaussian Auto-regressive models

The auto-regressivity of the recurrent neural network was used to model $q(z_i|z_{<i})$:

$$q_\eta(\mathbf{z}) = \prod q(z_i|z_{<i}) = \prod \mathcal{N}(z_i|\mu(z_{<i}), \sigma(z_{<i}))$$

Algorithm 4 Pseudocode for estimating the marginal KL using a density estimator for $q(\mathbf{z})$

- 1: Load trained variational model with prior $p(\mathbf{z})$, and posterior $q_\eta(\mathbf{z}|\mathbf{x})$.
 - 2: Initialize density \mathbf{t} model parameters ω randomly.
 - 3: **for** iter = 1 : max_iter **do**
 - 4: Update parameters ω by maximizing $\mathbb{E}_{p^*(\mathbf{x})} \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} [\log(\mathbf{t}(\mathbf{z}))]$
 - 5: **end for**
 - 6: marginal_kl = 0.0
 - 7: **for** i = 1 : num_z **do**
 - 8: sample \mathbf{x}_i from $p^*(\mathbf{x})$, sample \mathbf{z}_i from $q_\eta(\mathbf{z}|\mathbf{x}_i)$
 - 9: marginal_kl+ = $\log(\mathbf{t}(\mathbf{z}_i)) - \log p(\mathbf{z}_i)$
 - 10: **end for**
 - 11: marginal_kl = marginal_kl/num_z
-

F VGH++ loss function and pseudocode

In all the equations below, D_θ refers to the data discriminator, while C_ω refers to the code discriminator. They are both trained to minimize a cross entropy loss, like in the original GAN formulation.

The loss function for VGH is:

$$\mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} \left[-\lambda \|\mathbf{x} - \mathcal{G}_\theta(\mathbf{z})\|_1 + \log \frac{\mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z}))}{1 - \mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z}))} + \log \frac{C_\omega(\mathbf{z})}{1 - C_\omega(\mathbf{z})} \right]$$

In contrast, the loss function for VGH++ is:

$$\mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})} \left[-\lambda \|\mathbf{x} - \mathcal{G}_\theta(\mathbf{z})\|_1 + \log \frac{\mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z}))}{1 - \mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z}))} + \log \frac{C_\omega(\mathbf{z})}{1 - C_\omega(\mathbf{z})} \right] + \mathbb{E}_{p(\mathbf{z})} \log \frac{\mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z}))}{1 - \mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z}))}$$

The overall training procedure is summarized in Algorithm 5.

G The effect of the visible distribution in VAE training

Figure 22 visually shows the trade-off seen between using a Bernoulli or a QuantizedNormal distribution as the visible pixel distribution, $p(\mathbf{x}|\mathbf{z})$ in VAEs.

We now unpack the mathematical justification for why the Bernoulli distribution produces worse reconstructions. We will perform the analysis for a pixel x , but this straightforwardly extends to entire images. Assume a Bernoulli distribution with mean $\mu \in [0, 1]$. Then the Bernoulli loss is $x\mu - \log(1 + e^\mu)$. The gradient of the loss is $x - \sigma(\mu)$, where σ is the sigmoid

function. If we use a Gaussian distribution with mean m and standard deviation s , the gradient is $\frac{x-m}{s^2}$. We can see that the gradients of the two distributions have the same form, and using a Bernoulli distribution is equivalent to using a Gaussian distribution with variance 1. By setting the variance to 1, using a Bernoulli likelihoods spreads mass around each pixel and cannot specialize to produce good reconstructions.

H Tracking the variational lower bound - VEEGAN

In this section we compare the evidence lower bound training behavior, between VAEs and VEEGAN. VEEGAN introduces a new bound,

$$\text{KL}[p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})||p_\gamma(\mathbf{z}|\mathbf{x})p(\mathbf{x})] - \mathbb{E}[\log(p(\mathbf{z}))] + l_1(\mathbf{z}, F_\theta(\mathbf{x})) \quad (15)$$

where F_θ is the reconstructor network used to build the implicit $p(\mathbf{x}|\mathbf{z})$ and the KL divergence is estimated using the density ratio trick. The expected and desired behavior is that the bound increases as training progresses, however, as seen in Figure 23, variational hybrids do not solve one of the fundamental problems with adversarial model training, namely introducing a quantity to use to assess convergence.

I Real data evaluation metrics

A universal metric that can assess both overfitting, sample quality and sample diversity has not been found for generative models. Instead, multiple metrics which assess different aspects of a model have been proposed. To get a better overview of model performance, we use metrics which each capture a different aspect of training.

Inception score - sample quality and between class sample diversity: The most popular evaluation metric for implicit models is the Inception Score (Salimans et al., 2017). The Inception Score correlates with human sample evaluation and measures sample quality, between class sample diversity, but cannot capture withing class mode collapse (for example, the model could generate the same horse again and again, and the Inception Score will not penalize it) or overfitting. The Inception score uses the last layer logits of a ImageNet trained Inception network (Szegedy et al., 2016) to determine how classifiable a sample is. For generative models trained on CIFAR10, we complement our reporting by using a VGG style convolutional neural network, trained on CIFAR10, which obtained 5.5% error.

Multi-scale structural similarity (MS-SSIM): sample diversity: To measure sample diversity, we

Algorithm 5 Pseudocode for VGH++

- 1: Initialize parameters of generator θ , encoder η , discriminator ϕ and code discriminator ω randomly.
- 2: Let $\hat{\mathbf{z}} \sim q_\eta(\mathbf{z}|\mathbf{x})$ denote a sample from $q_\eta(\mathbf{z}|\mathbf{x})$ and $\hat{\mathbf{x}} = \mathcal{G}_\theta(\hat{\mathbf{z}})$ denote the ‘reconstruction’ of \mathbf{x} using $\hat{\mathbf{z}}$.
- 3: Let $R_{\mathcal{D}_\phi}(\mathbf{x}) = -\log \mathcal{D}_\phi(\mathbf{x}) + \log(1 - \mathcal{D}_\phi(\mathbf{x}))$
- 4: Let $R_{\mathcal{C}_\omega}(\mathbf{z}) = -\log \mathcal{C}_\omega(\mathbf{z}) + \log(1 - \mathcal{C}_\omega(\mathbf{z}))$
- 5: **for** iter = 1 : max_iter **do**
- 6: Update encoder η by minimizing

▷ data reconstruction and code generation loss

$$\mathbb{E}_{p^*(\mathbf{x})}\mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})}[\lambda\|\mathbf{x} - \mathcal{G}_\theta(\mathbf{z})\|_1] + R_{\mathcal{C}_\omega}(\mathbf{z}) \approx \mathbb{E}_{p^*(\mathbf{x})}[\lambda\|\mathbf{x} - \hat{\mathbf{x}}\|_1] + R_{\mathcal{C}_\omega}(\hat{\mathbf{z}}) \quad (13)$$

- 7: Update generator θ by minimizing

▷ data reconstruction and generation loss

$$\begin{aligned} & \mathbb{E}_{p^*(\mathbf{x})}\mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})}[\lambda\|\mathbf{x} - \mathcal{G}_\theta(\mathbf{z})\|_1 + R_{kl}(\mathbf{z})] + \mathbb{E}_{p(\mathbf{z})}[R_{\mathcal{D}_\phi}(\mathcal{G}_\theta(\mathbf{z}))] \\ & \approx \mathbb{E}_{p^*(\mathbf{x})}[\lambda\|\mathbf{x} - \hat{\mathbf{x}}\|_1 + R_{\mathcal{D}_\phi}(\hat{\mathbf{x}})] + \mathbb{E}_{p(\mathbf{z})}[R_{\mathcal{D}_\phi}(\mathcal{G}_\theta(\mathbf{z}))] \end{aligned}$$

- 8: Update discriminator ϕ by minimizing

▷ treat data as real, reconstructions and samples as fake

$$\begin{aligned} & \mathbb{E}_{p^*(\mathbf{x})}[-2\log \mathcal{D}_\phi(\mathbf{x}) - \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})}\log(1 - \mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z})))] + \mathbb{E}_{p(\mathbf{z})}[-\log(1 - \mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z})))] \\ & \approx \mathbb{E}_{p^*(\mathbf{x})}[-\log \mathcal{D}_\phi(\mathbf{x}) - \log(1 - \mathcal{D}_\phi(\hat{\mathbf{x}}))] + \mathbb{E}_{p(\mathbf{z})}[-\log(1 - \mathcal{D}_\phi(\mathcal{G}_\theta(\mathbf{z})))] \end{aligned}$$

- 9: Update code discriminator ω by minimizing

▷ treat $p(\mathbf{z})$ as real and codes from the encoder as fake

$$\mathbb{E}_{p^*(\mathbf{x})}\mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x})}[-\log(1 - \mathcal{C}_\omega(\mathbf{z}))] + \mathbb{E}_{p(\mathbf{z})}[-\log \mathcal{C}_\omega(\mathbf{z})] \approx \mathbb{E}_{p^*(\mathbf{x})}[-\log(1 - \mathcal{C}_\omega(\hat{\mathbf{z}}))] + \mathbb{E}_{p(\mathbf{z})}[-\log(\mathcal{C}_\omega(\mathbf{z}))] \quad (14)$$

10: **end for**

use 1.0 - MS-SSIM (Wang et al., 2003), an image similarity metric ranging between 0.0 (low similarity) and 1.0 (high similarity) that has been shown to correlate well with human judgment. The use of MS-SSIM for sample diversity was introduced by Odena et al. (2016), which used it compute in class sample similarity for conditional models, as between class variability can lead to ambiguous results. For models trained on the CelebA dataset, we can use this sample diversity metric, since the dataset only contains faces.

Independent Wasserstein critic - sample quality and overfitting: Danihelka et al. (2017) and Jiwoong Im et al. (2018) proposed training an independent Wasserstein GAN critic to distinguish between real data and generated samples. This metric has been shown to correlate with human evaluations (Jiwoong Im et al., 2018), and if the independent critic is trained on validation data, it can also be used to measure overfitting (Danihelka et al., 2017). All our reported results using the Independent Wasserstein Critic use a WGAN-GP model, trained to distinguish between the data validation set and model samples.

J AdversarialVB results

Results obtained using AdversarialVB are presented in Figure 24.

K Training details: hyperparameters and network architectures

For all our models, we kept a fixed learning rate throughout training. We note the difference with AGE, where the authors decayed the learning rate during training, and changed the loss coefficients during training²). The exact learning rate sweeps are defined in Table 3. We used the Adam optimizer (Kingma and Ba, 2015) with $\beta_1 = 0.5$ and $\beta_2 = 0.9$ and a batch size of 64 for all our experiments. We used batch normalization (Ioffe and Szegedy, 2015) for all our experiments. We trained all ColorMNIST models for 100000 iterations, and CelebA and CIFAR-10 models for 200000 iterations.

Scaling coefficients

We used the following sweeps for the models which have combined losses with different coefficients (for all our baselines, we took the sweep ranges from the original papers):

- WGAN-GP

²As per advice found here: <https://github.com/DmitryUlyanov/AGE/>

Network	VAE	DCGAN WGAN-GP	AAE	VGH VGH++	VEEGAN
Generator/Encoder	0.001, 0.0005, 0.005	0.0001, 0.0002, 0.0003	0.001, 0.0005, 0.005	0.0001, 0.0005	0.001, 0.0005
Discriminator	—	0.0001, 0.0002, 0.0003	—	0.0005	0.00005, 0.0001
Code discriminator	—	—	0.0005, 0.00005, 0.00001	0.0005	—

Table 3: Learning rate sweeps performed for each model.

- The gradient penalty of the discriminator loss function: 10.
- VGH++ and VGH
 - Data reconstruction loss for the encoder: sweep over 1, 5, 10, 50.
 - Data reconstruction loss for the generator: sweep over 1, 5, 10, 50.
 - Adversarial loss for the generator (coming from the data discriminator): 1.0.
 - Adversarial loss for the encoder (coming from the code discriminator): 1.0.

For Adversarial Autoencoders and VEEGAN, we followed the advice from the original paper and did not weight the different loss terms using coefficients.

Choice of loss functions

For VEEGAN, we used the l_1 loss as the code reconstruction loss. For VGH and VGH++ , we used l_1 as the data reconstruction loss and the classifier GAN loss for the data and code discriminator.

Updates

For the WGAN-GP experiments, we did 5 discriminator updates for generator update. All other models used the same number updates for model component (discriminator, generator, encoder, decoder).

Choice of latent prior

We use a univariate normal prior for all models.

K.1 Network architectures

For all our baselines, we used the same discriminator and generator architectures, and we controlled the number of latents for a fair comparison. For methods which needs an encoder such as VAEs, VEEGAN, VGH and VGH++ , the encoder is always set as a convolutional network, formed by transposing the generator (we do not use any activation function after the encoder). All discriminators use leaky units (Maas et al., 2013) with a slope of 0.2, and all generators used ReLUs. In all VAE results, unless otherwise specified, we used a Bernoulli visible distribution and a Gaussian latent posterior.

ColorMNIST

For all our models trained on ColorMNIST, we swept

over the latent sizes 10, 50 and 75. Tables 4 and 5 describe the discriminator and generator architectures respectively.

CelebA and CIFAR-10

The discriminator and generator architectures used for CelebA and CIFAR-10 were the same as the ones used by Gulrajani et al. (2017) for WGAN, using code at <https://github.com/martinarjovsky/WassersteinGAN/blob/master/models/dcgan.py>. Note that the WGAN-GP paper reports Inception Scores computed on a different architecture, using 101-Resnet blocks. For VEEGAN, we designed a code discriminator as defined in Table 6.

Code discriminator architectures

For a fair comparison between models, we used the same code discriminator architecture, where one is applicable. We tried both deeper convolutional architectures as well as shallow but bigger linear layers. We found the latter to work best and hence we used a 3 layer MLP with 1000 units each and leaky RELUs activations (Maas et al., 2013) with a slope of 0.2 as the code discriminator. Using 5000 units did not substantially improve results. This can be explained by the fact that too strong gradients from the code discriminator can effect the reconstruction ability of the encoder, and then more careful tuning of loss coefficient is needed. Perhaps optimization algorithms which are better suited for multi loss objective could help this issue.

Operation	Kernel	Strides	Feature maps
Convolution	5×5	2×2	8
Convolution	5×5	1×1	16
Convolution	5×5	2×2	32
Convolution	5×5	1×1	64
Convolution	5×5	2×2	64
Linear adv	—	—	2
Linear class	—	—	10

Table 4: ColorMNIST data discriminator architecture used for all models which require one. For DCGAN, we use dropout of 0.8 after the last convolutional layer. No other model uses dropout.

Operation	Kernel	Strides	Feature maps
Linear	—	—	3136
Transposed Convolution	5×5	2×2	64
Transposed Convolution	5×5	1×1	32
Transposed Convolution	5×5	2×2	3

Table 5: ColorMNIST generator architecture. This architecture was used for all compared models.

Operation	Kernel	Strides	Feature maps
Convolution	3×3	1×1	[512, 1024]
Convolution	2×2	1×1	[512, 1024]
Linear adv	—	—	2

Table 6: The joint discriminator head used for VEE-GAN. The input of this network is a vector concatenation of data and code features, each obtained by passing the data and codes through the same discriminator and code architectures used for the other models (excluding the classification layer).

L Reconstructions

We show reconstructions obtained using VGH++ and VAEs for the CelebA dataset in Figure 25 and on CIFAR-10 in Figure 26.

M Model samples for real data experiments

We show samples obtained on CelebA, CIFAR10 and ColorMNIST in Figures 27, 28 and 29, respectively.

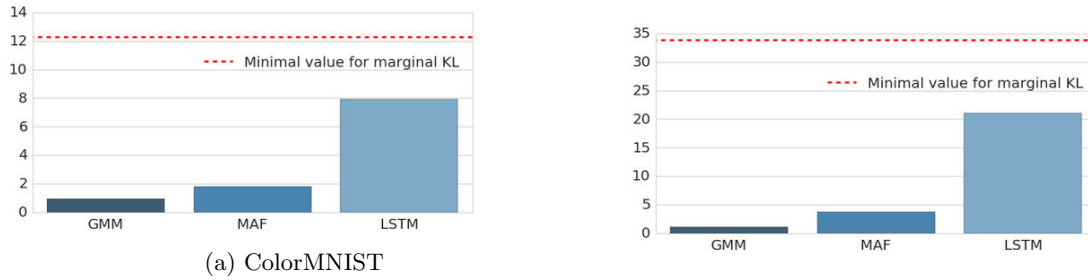


Figure 21: Estimating $KL[q_\eta(\mathbf{z})||p(\mathbf{z})]$ using different density models to estimate $q(z)$: Gaussian Mixture Models (GMM), Masked Auto-regressive Flow (MAF) and autoregressive models (LSTM). We plot the minimal value for the marginal KL - computed from Equation 5 - which allows us to conclude that all three density estimation approaches underestimate the true KL. LSTMs outperform the other models, showing that the autoregressivity of these models is necessary to model $q_\eta(\mathbf{z})$.

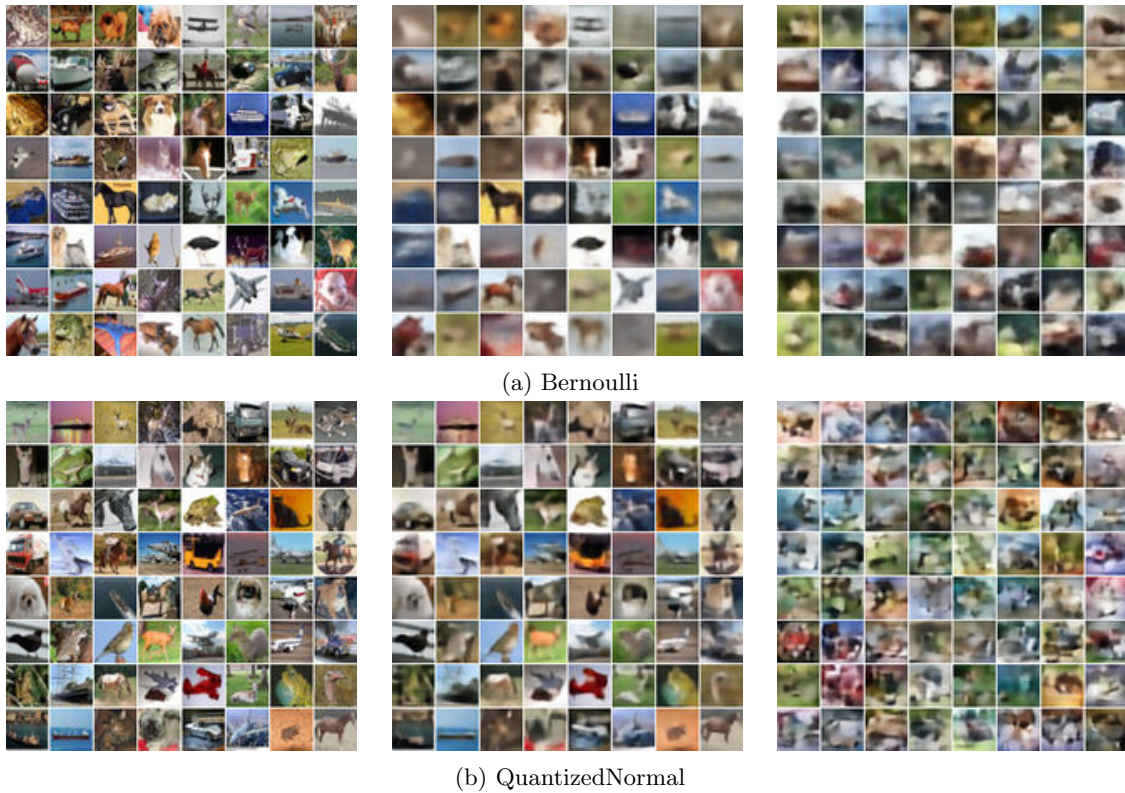


Figure 22: Comparisons of reconstructions and samples generated from VAEs using Bernoulli and Quantized Normal pixel distributions. The observed trade-off between reconstruction and sample quality is consistent throughout different hyperparameters. For the models displayed here, the difference can be seen in the different KL values obtained in the loss function used to train the models: 44.7 for the Bernoulli model, and 256.7 for the QuantizedNormal model.

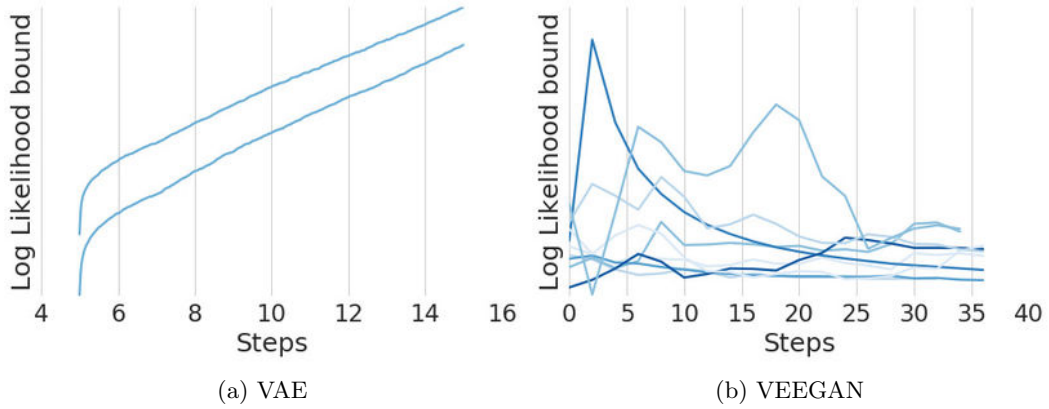


Figure 23: CIFAR-10 training curves of a standard VAE and VEEGAN across different hyperparameters. Results were obtained on . Since the bound of VEEGAN is not obtained on the observed data the numbers are not directly comparable. The aim of this plot is to show the **trend** of training, as we expect that as model training progresses, the likelihood increases. We see that for VEEGAN this is not the case, even though the models perform comparable with state of the art (Figure 8).

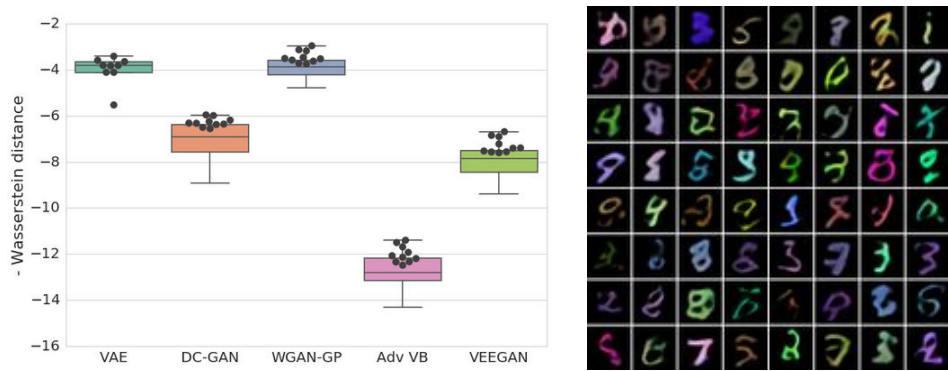
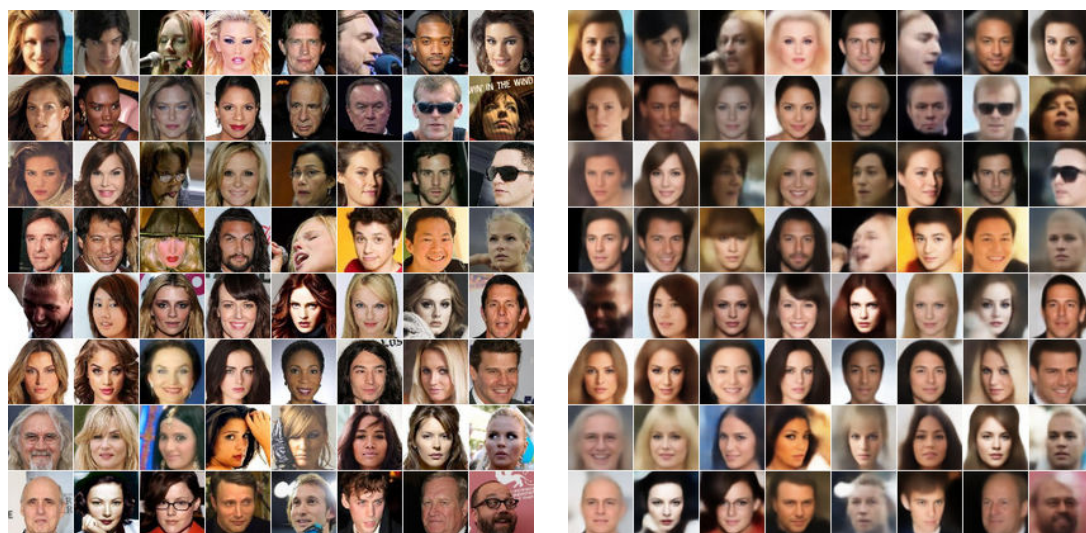


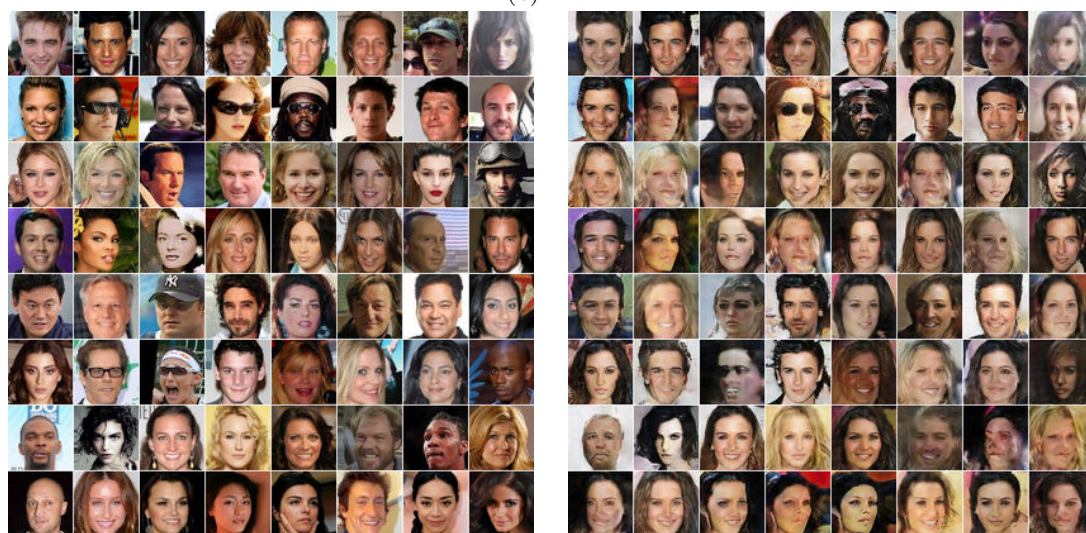
Figure 24: Results obtained using Adversarial VB, without adaptive contrast. The model was not able to match $q_{\eta}(\mathbf{z})$ and $p(\mathbf{z})$, and this results in an independent critic being able to easily distinguish samples from data (left). This can be seen visually on the right, as the digits do not appear well defined for a large number of samples.



(a) VAE

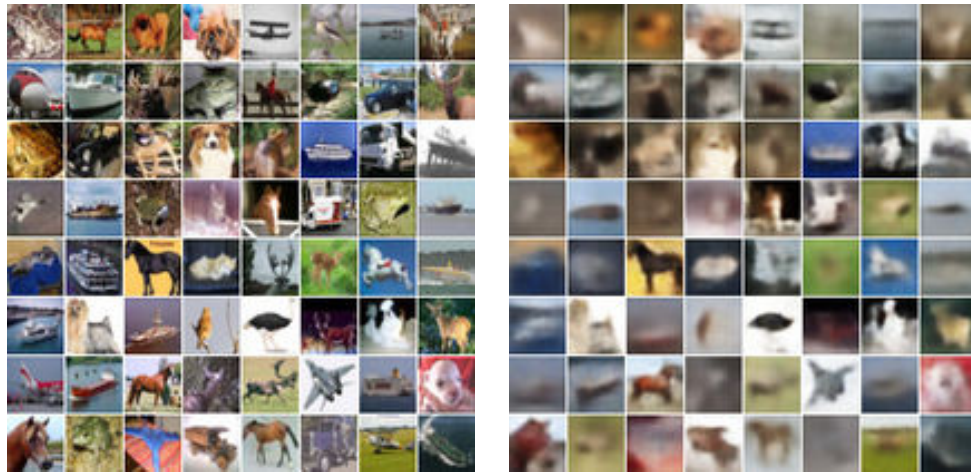


(b) AAE

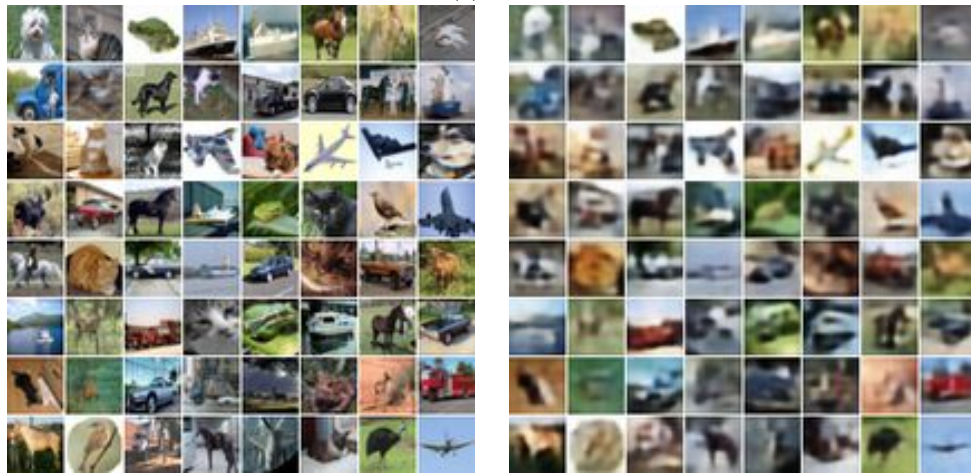


(c) VGH++

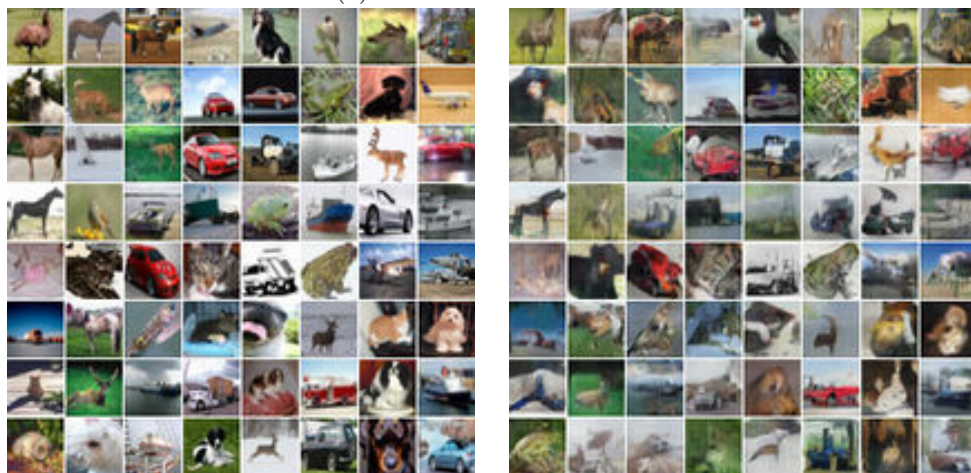
Figure 25: Training reconstructions obtained using a standard VAE, Adversarial Autoencoders and VGH++ on CelebA. Left is the data and right are reconstructions.



(a) VAE



(b) Adversarial Autoencoders



(c) VGH++

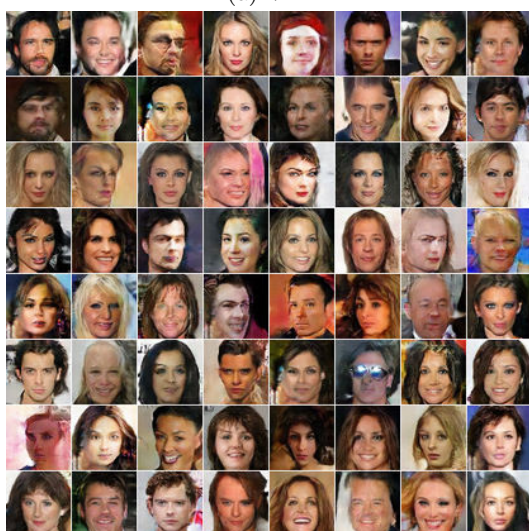
Figure 26: Training reconstructions obtained using a standard VAE, Adversarial Autoencoders and VGH++ on CIFAR-10. Left is the data and right are reconstructions.



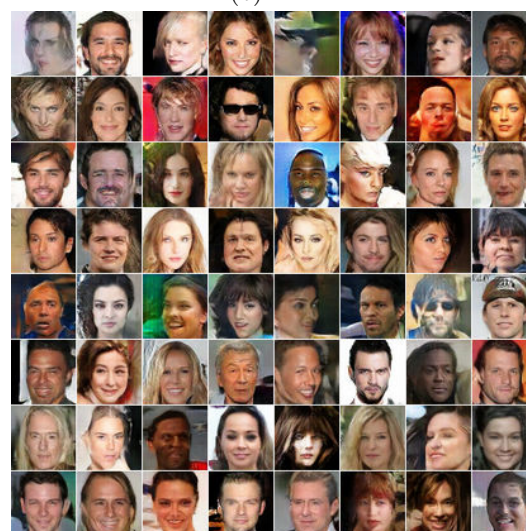
(a) VAE



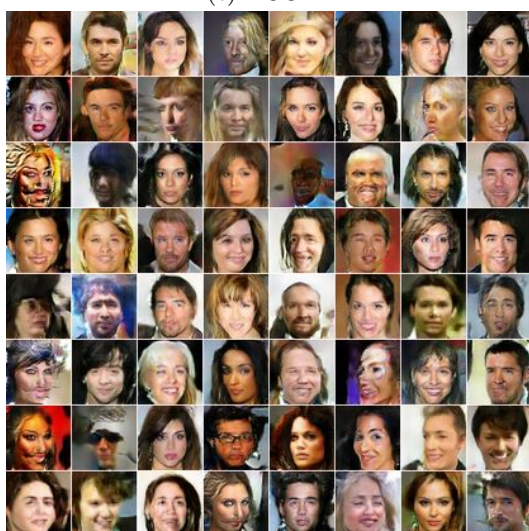
(b) AAE



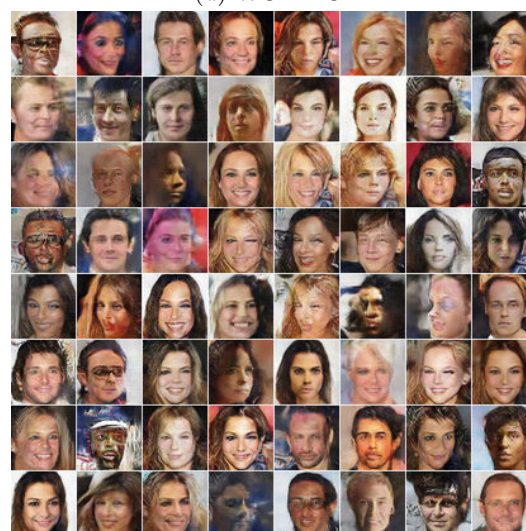
(c) DCGAN



(d) WGAN-GP



(e) VEEGAN



(f) VGH++

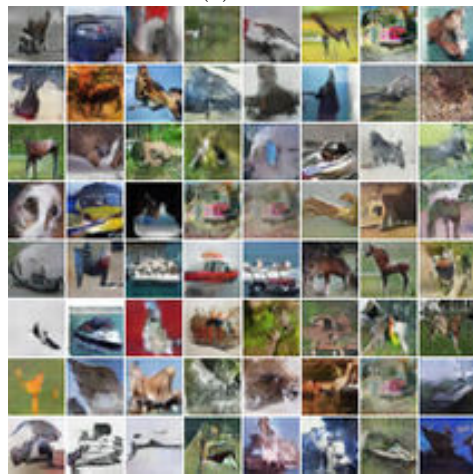
Figure 27: CelebA samples.



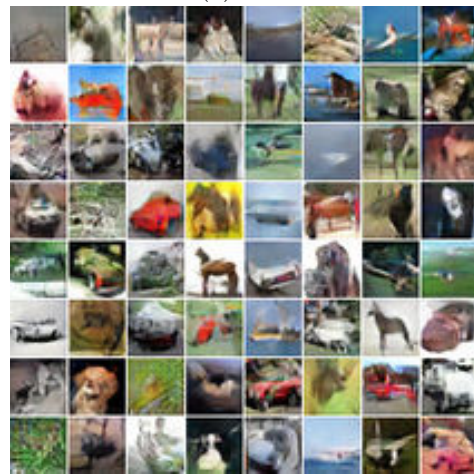
(a) VAE



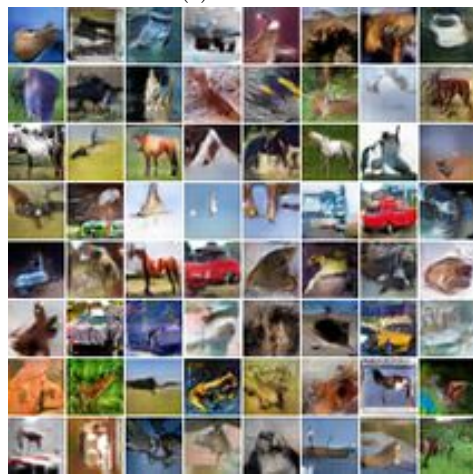
(b) AAE



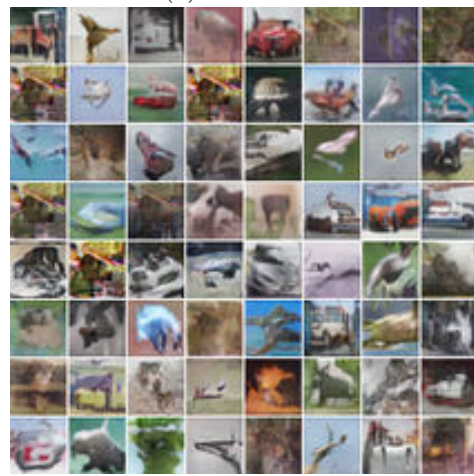
(c) DCGAN



(d) WGAN-GP



(e) VEEGAN



(f) VGH++

Figure 28: CIFAR10 samples.



(a) VAE



(b) AAE



(c) DCGAN



(d) WGAN-GP



(e) VEEGAN



(f) VGH++

Figure 29: ColorMNIST samples.