# *Jigsaw*-VAE: Towards Balancing Features in Variational Autoencoders

**Saeid Asgari Taghanaki**
Autodesk AI Research
Toronto, Canada

**Mohammad Havaei**[*]
Imagia
Montreal, Canada

**Alex Lamb**[*]
MILA
Montreal, Canada

**Aditya Sanghi**
Autodesk AI Research
Toronto, Canada

**Ara Danielyan**
Autodesk AI Research
Toronto, Canada

**Tonya Custis**
Autodesk AI Research
San Francisco, USA

## Abstract

The latent variables learned by VAEs have seen considerable interest as an unsupervised way of extracting features, which can then be used for downstream tasks. There is a growing interest in the question of whether features learned on one environment will generalize across different environments. We demonstrate here that VAE latent variables often focus on some factors of variation at the expense of others - in this case we refer to the features as "imbalanced". Feature imbalance leads to poor generalization when the latent variables are used in an environment where the presence of features changes. Similarly, latent variables trained with imbalanced features induce the VAE to generate less diverse (i.e. biased towards dominant features) samples. To address this, we propose a regularization scheme for VAEs, which we show substantially addresses the feature imbalance problem. We also introduce a simple metric to measure the balance of features in generated images.

## 1 Introduction

Variational autoencoders (VAEs) have certain appealing properties compared to generative adversarial networks (GANs), such as stable training [37], interpretable inference [5], and calculating data likelihood [18]. As such, they remain worthy for examination and improvement. Additionally, recent VAEs [30] have shown a great potential for generating competitive high quality images compared to the state-of-the-art GANs.

However, there are a few identified challenges with regard to VAEs: balancing the two terms in the VAE loss function (i.e. the log likelihood and the Kullback–Leibler (KL) divergence terms) is not trivial. Sacrificing the former causes sub-optimal reconstruction performance while neglecting the latter results in poor sampling (several latent variables might be ignored/over-pruned [39]). To tackle this, several works have explored weighting the terms in a more systematic way than the original VAE e.g. by annealing the contribution weight [3, 34] or by learning it [1, 8]. Additionally, even after convergence, a matched latent prior with a learned aggregated posterior distribution is not guaranteed, which can be due to the choice of an overly simplistic prior distribution [19, 8]. A simple prior and/or improperly weighted evidence lower bound (ELBO) terms can lead to less diverse generated samples [38, 14] i.e. neglecting the minor (sub-) clusters of input samples/features.

In this paper, we study VAEs from a slightly different perspective. We examine whether optimized VAEs via the ELBO tend to ignore sporadic features of input data since the models are able to optimize
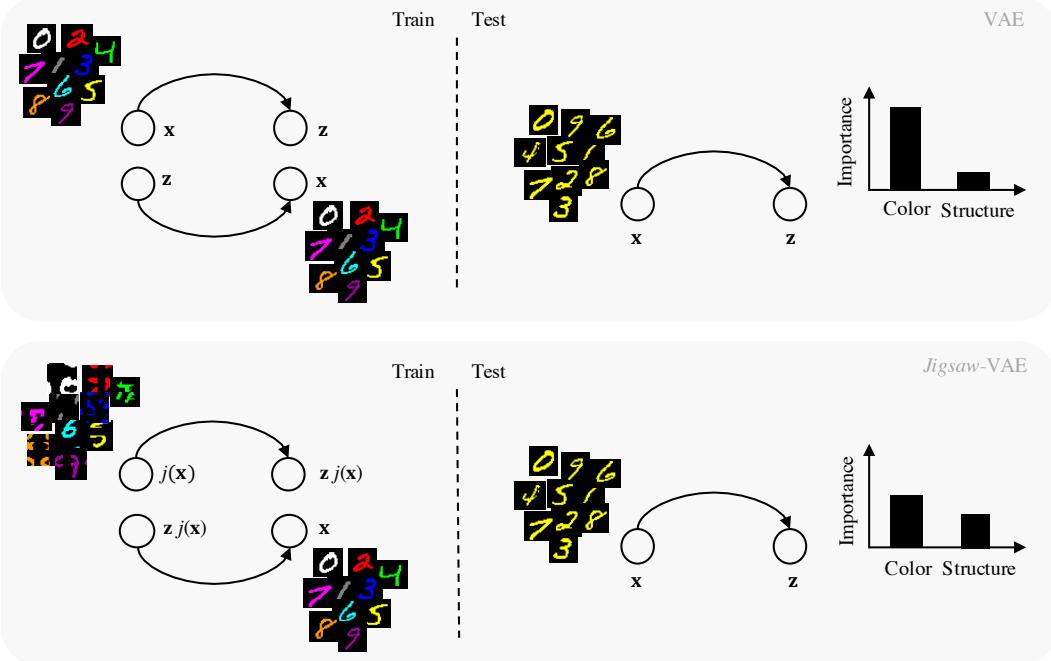
---

[*]Equal contribution

Figure 1: Colored MNIST example of balancing features : comparison between the VAE and our proposed *Jigsaw*-VAE. Note that the permutation function $J$ is applied on each individual sample separately.

the ELBO faster by comfortably favoring the optimization towards more frequent, dominant features. Handling *feature imbalance* (a.k.a representation bias [23, 36] in general) in variational models, particularly VAEs, seems to be more complex than supervised and non-variational models [41, 12, 21, 16, 11, 25, 7, 22] where they have access to the prior information/distribution of features or they specify a target feature (e.g. color) to balance. However, it is not trivial to find problematic features in a large dataset to target while performing a generative task or representation learning. Similarly, the classic problem of class imbalance in supervised classification models can be handled with several different approaches such as by simply adding penalty terms/weights to the model's loss function to enforce the model for a fair, balanced prediction. However, when it comes to generative models and representing learning, there can be many large and small clusters of samples in training data (usually without fine-grained cluster labels). Here, the question to be answered is: *what does imbalance mean in generative models and representation learning and how it can be handled?* In other words, for the same example above, how do we enforce a generative model to not ignore a non-dominant cluster while generating random samples? Note that even enforcing the model to handle the high level imbalance of "female" vs "male", the same issue can be exacerbated by extending to sub-clusters of each of these categories e.g. females with and without eyeglasses, blond or black hair.

Therefore, besides the main challenges of VAEs discussed above, feature imbalance might also be a reason for capturing poor latent information which results in generating less diverse samples as well as under-performing in downstream tasks which are conditioned on latent features. For instance, if we assume a data distribution is defined by two main features of "color" and "structure", we would ideally want a model to learn "equal" contribution of the two features. However, this does not always happen in practice, which might result in a drastic failure in downstream tasks. As an example, a simple case scenario has been visualized in Figure 1 (top) where modification of dominant feature (color) results in a feature-imbalanced latent representation. As a remedy, we apply feature permutation (Figure 1 (bottom)) to the VAE to reduce reliance on a single or a few dominant feature(s).

Feature permutation has previously been used as a metric to calculate feature importance in Random Forests [4]. A feature is "important" if permuting its values decreases the model performance, indicating that the model had relied on that feature for its prediction. Based on this idea, Fisher et al. [10] proposed a model-agnostic version of the feature importance. Recently, permutation

based techniques, particularly *Jigsaw* puzzle approaches, have been successfully applied to (semi-) supervised deep models [28, 33, 29] with the goal of capturing rich contextual information. However, their aptitude for addressing the feature imbalance in the latent space of VAEs and how these latent spaces subsequently affect the downstream tasks has not yet been studied. In this paper, we focus on *whether permutation-based regularization is capable of reducing the adverse effect of feature imbalance in VAEs*.

To solve a jigsaw puzzle, a VAE has to focus on how different input variables can be integrated (sorted) to build a plausible structure. In other words, instead of reconstructing the "big picture" by focusing on dominant features which might be a simpler task to optimize the ELBO, the model focuses more on smaller areas that include local structural features. This prevents the model to choose the easy optimization path of learning *only* dominant variables.

In this paper, we make the following contributions:

- We inspect the feature imbalance issue in the context of VAE. We study whether the VAE learns balanced features given feature-imbalanced input data, which is more often the case in practice. Training a VAE with feature-balanced data is the ideal way of achieving this goal. However, collecting feature-balanced data requires great effort.

- We propose a simple metric called *Feature Presence Metric* to measure the balance of features in generated images.

- We propose a simple yet effective feature permutation-based technique to systematically enforce the VAE to learn balanced features.

- On the line of research of improving the VAE using more flexible priors, we introduce a prior which is a mixture of a uniform permutation distribution and a Gaussian distribution.

## 2 Method

In vanilla VAE, the prior distribution $p(\mathbf{z})$ is defined on the latent representation $\mathbf{z} \in \mathbb{R}^D$ and is usually set to an isotropic Gaussian distribution $\mathcal{N}(0, I)$ where $D$ is the latent dimension. The posterior distribution is defined as $p_\theta(\mathbf{z}|\mathbf{x}) \propto p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})$. Then a parameterized distribution over $\mathbf{x}$ is defined as $p_\theta(\mathbf{x}|\mathbf{z})$ and is modeled as a generative network in the context of neural network such that $\theta$ becomes the weights of the network. Similar to the generative network, a neural network is used to approximate distribution $q$ conditioned on observation $\mathbf{x}$, called inference network $q_\phi(\mathbf{z}|\mathbf{x})$ with variational parameter $\phi$ which is also weights of the neural network. Using re-parameterization trick [18] back-propagation is applied on the parameter $\phi$ considering $\mathbf{z}$ as a function of noise and typically mean and variance of the Gaussian learned by a decoder network. Therefore, the objective of VAE is to maximize the following variational lower bound with respect to the parameters $\theta$ and $\phi$,

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log p_\theta(\mathbf{x}|\mathbf{z}) \right] - \mathbb{KL} \left( q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}) \right) \quad (1)$$

Let $q_\phi(\mathbf{z}|j(\mathbf{x}))$ be a conditional normal distribution such that

$$q_\phi(\mathbf{z}|j(\mathbf{x})) = \mathcal{N}(\mathbf{z}|\mu_\phi(j(\mathbf{x})), \sigma_\phi(j(\mathbf{x}))). \quad (2)$$

with $\mu_\phi(j(\mathbf{x}))$ and $\sigma_\phi(j(\mathbf{x}))$ as non-linear functions and $j$ as a permutation function. If we consider $p(j(\mathbf{x})|\mathbf{x}) = \mathcal{U}_{[0,b]}(\mathbf{x})$ to be a permutation model (Jigsaw in our case) of $\mathbf{x}$ where $b$ is the number of different permutations in the uniform distribution, then

$$q_\phi^j(\mathbf{z}|\mathbf{x}) = \mathbb{E}_{p(j(\mathbf{x})|\mathbf{x})} \left[ q_\phi(\mathbf{z}|j(\mathbf{x})) \right] = \int_{j(\mathbf{x})} q_\phi(\mathbf{z}|j(\mathbf{x}))p(j(\mathbf{x})|\mathbf{x})dj(\mathbf{x}) \quad (3)$$

is a mixture of normal and uniform distributions (i.e. each time we sample $j(\mathbf{x}) \sim p(j(\mathbf{x})|\mathbf{x})$ and feed into $q(\mathbf{z}|j(\mathbf{x}))$ we get different distributions). In practice, the inference network will learn which distribution is needed for a given $j(\mathbf{x})$.

Similar to [15], if we consider $q_\phi^j(\mathbf{z}|\mathbf{x}))$ as the approximate distribution, the variational lower bound can be written as

$$\mathcal{L}_{Jigsaw\text{-}VAE} = \mathbb{E}_{q_\phi^j(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|j(\mathbf{x}))} \right]. \tag{4}$$

Compared to $q_\phi(\mathbf{z}|\mathbf{x})$ in VAE, $q_\phi^j(\mathbf{z}|\mathbf{x})$ in *Jigsaw*-VAE has the capacity to cover a broader class of distributions, thus creating more sub-spaces to decode diverse samples without intermixing sub-spaces. The comparison of the *Jigsaw*-VAE with the VAE and and other models which use complex priors such as VampPrior VAE [38] is depicted in Figure 2. As can be seen, the VampPrior VAE builds a hierarchical variational inference module by explicitly adding auxiliary latent vectors ($\mathbf{z}_1$, $\mathbf{z}_2$) while our *Jigsaw*-VAE has a single latent vector which caries information about both $\mathbf{z}$ and $\mathbf{z}' = \mathbf{z}(j(\mathbf{x}))$.
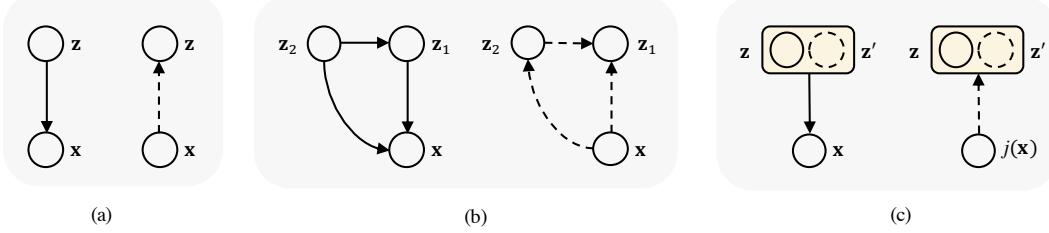


(a)       (b)       (c)

Figure 2: Stochastical dependencies. (a) VAE [18], (b) VampPrior VAE [38], (c) *Jigsaw*-VAE where $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})$ and $\mathbf{z}' \sim q(\mathbf{z}|j(\mathbf{x}))$. Solid and dashed arrows show generative and inference steps, respectively.

Compared to mixture of Gaussians prior and VampPrior, *Jigsaw*-VAE does not have any extra parameters beyond that required by the VAE to be learned. In the case of mixture of Gaussians, $p(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} \mathcal{N}\left(\mu_k, \mathrm{diag}\left(\sigma_k^2\right)\right)$ where $\mu_k \in \mathbb{R}^M, \sigma_k \in \mathbb{R}^M$ are learned parameters.

The permutation function we use as the input stochastic layer is depicted in Figure 3. The function can be defined over spatial windows $S$ with size $(\delta n, \delta m)$, color channels, or both.
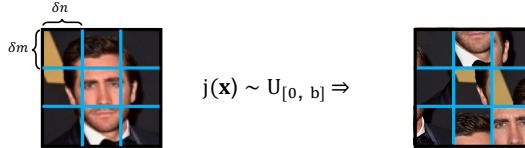


Figure 3: Sampling permutations form a uniform distribution. $(\delta n, \delta m)$ show the window size $S$. In our experiments, we set $\delta n = \delta m = \frac{I}{4}$ where $I$ corresponds to size of a side of the square input image. Therefore the range of the permutation function $j$ will be $[0, b = 16!]$ where $b$ is the total number of permutations.

## 3   Related Work

**Smoothing the conflict between the ELBO terms.** Ideally, we want every data sample to be (a) encoded into a latent space such that it can be accurately reconstructed and (b) close enough to the origin that we are likely to hit that point in latent space by sampling $\mathbf{z} \sim \mathcal{N}(0, I)$. Note that (a) and (b) are somewhat in conflict [14, 35]: enforcing the latent points to be far from each other makes (a) easier but harms (b). This conflict can be softened either by defining more complex priors than the Gaussian as done in VampPrior VAE [38] and denoising VAE [15], or by modifying the VAE objective [14]. Our proposed *Jigsaw*-VAE adds an extra stochastic sampling layer to the beginning of the encoder which makes the latent sampling space a mixture of distributions. However, in comparison with works [38, 15] which use more flexible priors than the standard Gaussian, our method does not need to modify the network architecture, nor does it need extra training with sudo inputs as done by Tomczak et al., [38]. In contrast to methods that modify the ELBO by introducing

a weighting coefficient for the KL term [14, 3, 34, 1, 8], our method does not require any balancing factor to determine the contribution of the likelihood and the KL term in the overall optimization.

**Stochastic layer in the bottom of encoder.** Compared to denoising VAE (*d*-VAE) [15] which adds an extra stochastic noise layer to the VAE, we use a jigsaw permutation function based layer to systematically enforce the capturing of structural features by the VAE. Our jigsaw approach encodes the prior assumptions into the model such that the structure of the world (relative position of parts) is more crucial to get right than the pixel-level appearance which it is done by *d*-VAE.

**Enforced information placement in latent vector.** The work done by Chen et al., [6] is related to our method in the sense that they have also induced the latent space to carry some desired information. To be specific, however, we make the decoding distribution of our method to be incapable of modeling information that the Jigsaw-VAE's encoder captures (i.e. thanks to the permutation function it codes the local spatial relationships). The vector quantized VAE [30] also filters the information in the latent space such that only the most discriminative features are passed through the decoder. Louizos et al., [26] have proposed a penalty term based on maximum mean discrepancy to obtain purged latent information with removed unwanted sources of variation. However, we prefer not to remove infrequent useful features. Instead, our method balances and recovers infrequent features in the latent space. Our work also differs from disentanglement approaches [27, 17]. Although they learn independent factors of variations, they do not guarantee that the less frequent factors will not be ignored or frequent features will be dominated.

In contrast to all of the above, to the best of our knowledge, our work is the first to study the feature imbalance concept in VAEs and its effect on downstream tasks. It is also the first work to explore the effect of a *permutation*-based stochastic layer, placed at the bottom of the VAE's encoder.

## 4 Experiments

In the experiments we aim at verifying empirically whether the *Jigsaw* approach (a) helps the VAE to learn a representation that can preserve rare features after convergence using ELBO in generative tasks, and (b) whether it helps learning balanced features for downstream tasks such as clustering.

Therefore we group the experiments into two main categories. To answer (a), in subsection 4.2, using our proposed metric (i.e. feature presence metric: subsection 4.1) we inspect the presence of targeted features in a generative task; to answer (b), in subsection 4.3, we measure the discriminate performance of the latent information in a feature-biased clustering scenario.

### 4.1 Feature Presence Metric

Various metrics have been proposed to evaluate the performance of VAEs. Among them, reconstruction error, negative log likelihood, Inception score (IS) [31] and Frechet Inception distance (FID) [13] (or a collection of these) are generally used. IS and FID are both based on a model trained on Imagenet [9] and as reported before, they do not properly capture fidelity or diversity [2] (e.g. unrealistic images can obtain high IS and FID scores). We argue that in addition to the above metrics (regardless of their limitations), measuring whether features/variables of the train set are preserved after a VAE is converged is an important factor towards both robust representation learning and diversity in generated samples. To this end, we propose the feature presence metric (FPM) which can be calculated over randomly generated images of a generative model. FPM measures how well a particular feature/variable is perpetuated in a set of generated images.

For a target feature $f$, FPM is calculated as:

$$FPM = \left| (\frac{N_{gf}}{N_g}) - (\frac{N_{tf}}{N_t}) \right| \times 100 \tag{5}$$

where $N_{gf}$ and $N_{tf}$ are the number of randomly generated and train (real) images which contain the target feature, respectively. $N_g$ and $N_t$ are the total number of the generated and train images, respectively. To calculate $N_{gf}$ we train a classifier which detects whether the targeted feature is present in a sample.

## 4.2 Feature Inspection

In this section, we use the CelebFaces Attributes (CelebA) dataset [24] to train the models and analyze the diversity of the images randomly generated by each method. We compare our *Jigsaw*-VAE to four other VAE methods: the vanilla VAE [18], *d*-VAE [15] $\beta$-VAE [14], and the VampPrior VAE [38]. Inspired from the Mixup approach [40] which was applied as a data augmentation in supervised tasks, we train the vanilla VAE with mixed inputs (*Mixup*-VAE) as another baseline. Particularly, we add stochasticity to the input of the VAE by mixing training samples using the linear interpolation mentioned in [40]. We then set the ground truth for the VAE to be one of the mixed samples which has the max contribution in the input mixed-up sample.

To evaluate the methods' generative performance in terms of FPM, we train MobileNet-V2 [32] to predict the attributes of each randomly generated image. We choose a diverse set of features including both under- and over-represented ones i.e. features that are present in only $\sim 2\%$ to $\sim 58\%$ of the training data (Table 1). We generate 5000 random samples per method and report the frequency of the targeted attributes present in the generated images.

In Table 1, we report the performance of each method on imbalanced features i.e. we measure whether each method is capable of preserving the balance of features. Ideally, FPM values for each feature should be zero. As can be seen the *Jigsaw*-VAE achieves favorable MSE and FPM values. $\beta$-VAE and VampPrior VAE obtain better values for some of the under-represented attributes as well as reconstruction error. However, looking at Figure 4, the random images generated via $\beta$-VAE and VampPrior VAE seem less realistic than our proposed *Jigsaw*-VAE. After applying our Jigsaw approach to $\beta$-VAE (i.e. *Jigsaw*-$\beta$-VAE), the generated images look more realistic compared to $\beta$-VAE and achieve better FPM values. Although the VampPrior VAE leverages more complex latent calculations, it does not scale to high dimensional images; as shown in Figure 4-(f). Considering FPM values and whether the randomly generated images look realistic, jigsaw based approaches outperform others. This experiment showed that the $\beta$-VAE is able to preserve under-represented features to some extent, but it generates unrealistic images. However, applying the jigsaw approach to both VAE and $\beta$-VAE leads to preserved underrepresented features as well as generating more realistic images.

Table 1: Reconstruction error (i.e. MSE) and FPM values (lower is better) for features from celebA dataset. Numbers in parenthesis below each feature name are percentage of the data that has the feature. In each column, two best values are highlighted.

| Method | MSE | Female (58.06) | Eyeglasses (6.46) | Bald (2.28) | Beard (16.58) | Smiling (47.97) | Gray (4.24) | AVG |
|---|---|---|---|---|---|---|---|---|
| VAE [18] | 0.0272 | 6.30 | 6.24 | 2.04 | 14.22 | 14.31 | 2.76 | 9.762 |
| *d*-VAE [15] | 0.0271 | 17.52 | 6.36 | 2.10 | 13.34 | 17.00 | 2.18 | 9.748 |
| $\beta$-VAE [14] | 0.0265 | 13.60 | 5.38 | 1.22 | 10.72 | 11.33 | 3.98 | 7.705 |
| VampPrior VAE [38] | 0.0107 | 30.12 | 4.80 | 2.10 | 14.44 | 5.97 | 4.20 | 10.270 |
| *Mixup*-VAE | 0.0310 | 18.70 | 6.34 | 1.98 | 13.96 | 14.17 | 3.06 | 9.702 |
| *Jigsaw*-VAE (ours) | 0.0273 | 11.26 | 6.00 | 1.60 | 14.10 | 10.27 | 0.74 | 7.328 |
| *Jigsaw*-$\beta$-VAE (ours) | 0.0287 | 13.00 | 5.38 | 0.86 | 9.96 | 11.81 | 2.46 | 7.252 |

We next examine the proposed method's ability to produce smooth interpolation (i.e. gradual transition). In Figure 5, we show an interpolation sample between with eyeglasses and without eyeglasses. As shown in the third row of Figure 5, *Jigsaw*-VAE gradually removes the eyeglasses from the image, moving from left to right. However, in the VAE interpolation outputs as shown in the first row, the eyeglasses feature is suddenly disappeared after the second sample from left. The VampPrior VAE (second row), performs the worse by completely ignoring the eyeglasses attribute. $\beta$-VAE (the second row) encourages gradual transition compared to the VAE, however not as much as *Jigsaw*-VAE does.

## 4.3 Quality of the Learned Latent Vectors in Clustering Context.

Compared to supervised models, their unsupervised counterparts are more susceptible to strong feature/distribution biases in the input data. That is, a dominant feature such as color or texture might
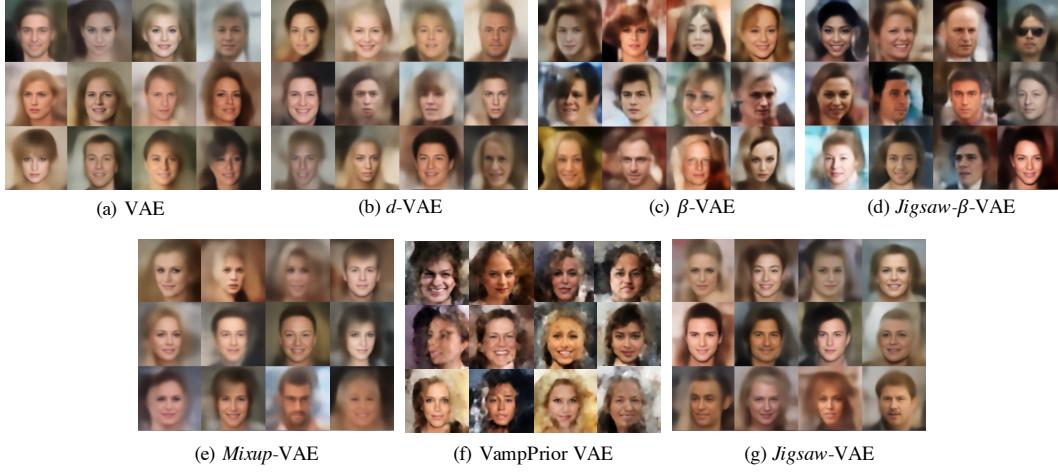
(a) VAE  (b) *d*-VAE  (c) *β*-VAE  (d) *Jigsaw-β*-VAE

(e) *Mixup*-VAE  (f) VampPrior VAE  (g) *Jigsaw*-VAE

Figure 4: Randomly generated samples of size $128 \times 128$ with different methods; zoom for details.



Figure 5: Transition from with to without eyeglasses. Rows from top to bottom correspond to VAE, VampPrior VAE, *β*-VAE, and *Jigsaw*-VAE, respectively.

confuse an unsupervised model to consider two different objects to be the same which have the same color/texture. In this section, in order to study the quality of the latent vectors in the VAE, we leverage the truncated Gaussian-mixture VAE [42]. We measure the clustering performance under an intended biased experiment; using colored MNIST dataset, we use a normal train set, but test with a biased set.

**Colored MNIST.** For this experiment, we colorize each digit class in the MNIST [20] train set with a specific color (Figure 1). To evaluate whether each method is able to capture more features besides the dominant *color* feature, we use a single color code for all the test images.

For this experiment, we compare the proposed *Jigsaw*-VAE to four other methods. Similar to the generative task experiments in subsection 4.2, we also train the VAE with the *Mixup* approach. In Table 2, we report the results of clustering experiments. These results indicate that although applying *β*-VAE and *d*-VAE improve generative task (subsection 4.2), they do not improve the performance on the normal colored test set. However, the *Mixup*-VAE and the proposed *Jigsaw*-VAE improve the normalized mutual information (NMI) [35] score from $0.8872$ to $0.9022$ and $0.9473$, respectively. However, when we test the methods with intended color bias (Single-color results in Table 2), all the methods except for the *Jigsaw*-VAE fail. This suggests that *Jigsaw*-VAE is able to balance shape and color features, thus not completely failing when the test set is highly biased towards the color feature. For this experiment, the stochastic *Jigsaw* permutation function runs over both spatial windows and RGB channels.

For this experiment, we adopt the truncated Gaussian-mixture VAE clustering approach [42] and modify it according to the different approaches used in the previous experiments i.e. VAE, *d*-VAE,

Table 2: Colored MNIST clustering results. Multi-color refers to test set which follows the roles of the train set i.e. each digit is assigned the same specific color as the train set. Single-color refers to the experiment where we assign a single color to all digits in the test set. Higher values better.

| | NMI | |
|---|---|---|
| Methods | Multi-color | Single-color (median) |
| VAE [18] | 0.8872 | 0.176 |
| d-VAE [15] | 0.8713 | 0.041 |
| β-VAE [14] | 0.8766 | 7.8E-06 |
| Mixup-VAE | 0.9022 | 0.114 |
| Jigsaw-VAE (ours) | 0.9473 | 0.379 |

*Mixup*-VAE, β-VAE, and *Jigsaw*-VAE. Ideally, we want the clusters created by each method to be based on both color and structural information. If a model assigns a correct cluster label to an input, the decoded information of the input should contain both correct color and shape features. In Figure 6, we visualize the reconstructed digits using clustering models trained with colored MNIST. When the test set color codes are the same as the training digit colors, cluster assignments are fairly accurate (Table 2) thus reconstructions are reasonable for all the methods (Figure 6 left). However, when all the digits in the test set take only one color code from train set, wrong clusters are assigned which results in wrong decodings (Figure 6 right). As can be seen, only *Jigsaw*-VAE is able to *recover* the actual color codes that the model was trained with. Since color "yellow" is the actual color for digit "5", non-Jigsaw methods tend to wrongly reconstruct 5-like images (second row, right) even though the input is digit "6". In another case (third row, right), non-Jigsaw approaches reconstruct a number form a closest color code to yellow, which in our case is "8" with color "orange". The results indicate that in the case of non-jigsaw VAE, color is learned as the dominant feature while the structural features are undermined. Where as the *Jigsaw*-vae, balances the structural and color features.
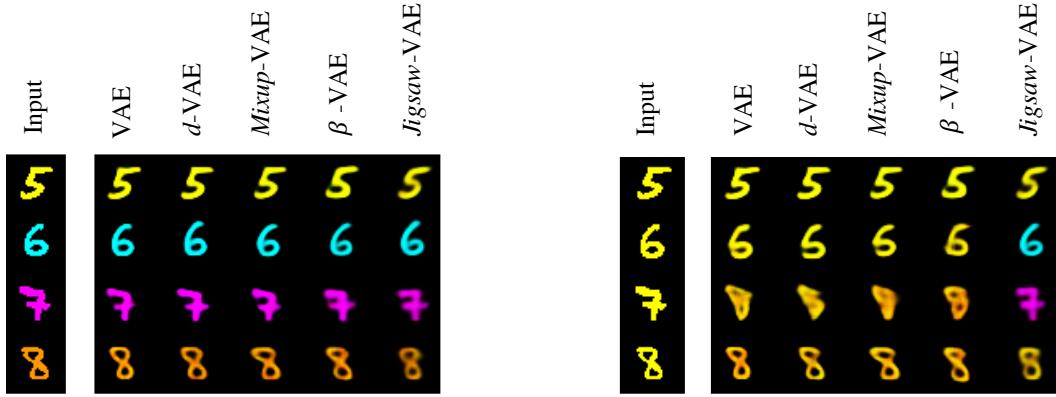


Figure 6: Colored MNIST reconstructed samples from decoding part of the clustering model [42]. Left: The input images (first column) have the same color codes as the train set distribution. Right: Input color codes are all mapped to a single code (here to yellow which is originally the color code for digit 5 in the train set) from train distribution.

## 5    Conclusion

We introduced *Jigsaw*-VAE, a new VAE capable of learning balanced features that can be used in both generative and other downstream tasks such as clustering. Our method exploits a mixture prior and a stochastic permutation layer. The mixture prior helps smoothing the conflict between the likelihood and KL terms while the permutation layer enforces the VAE to learn spatial relationships between the tiles (object parts). We empirically showed this prevents the VAE from ignoring infrequent spatial contexts (i.e. features). We showed that input permutation reduces the effect of strong damaging feature-bias in clustering models where we use VAEs as backbone. The *Jigsaw*-VAE also showed

smooth transitions during interpolation which is useful for preventing abrupt changes while modifying factors of a sample based on another sample.

Besides measuring reconstruction error, sampling quality and diversity, and how realistic the generated images are, we proposed to measure the presence of train data features in the generated images using our proposed metric, FPM. Among competing methods, $\beta$-VAE showed competitive performance on FPM and MSE metrics for the generative task (celebA experiments), however, it generates less realistic images compared to our *Jigsaw*-VAE. Additionally, $\beta$-VAE showed sub-optimal performance for the clustering task (MNIST experiments) when test set is biased. Therefore, considering both generative and clustering tasks, our *Jigsaw*-VAE outperformed competing methods.

# References

[1] Andrea Asperti and Matteo Trentin. Balancing reconstruction error and kullback-leibler divergence in variational autoencoders. *arXiv preprint arXiv:2002.07514*, 2020.

[2] Shane Barratt and Rishi Sharma. A note on the inception score. *arXiv preprint arXiv:1801.01973*, 2018.

[3] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.

[4] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.

[5] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. *arXiv preprint arXiv:1609.07093*, 2016.

[6] Xi Chen, Diederik P Kingma, Tim Salimans, Yan Duan, Prafulla Dhariwal, John Schulman, Ilya Sutskever, and Pieter Abbeel. Variational lossy autoencoder. *International Conference on Learning Representations*, 2017.

[7] Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. Don't take the easy way out: Ensemble based methods for avoiding known dataset biases. *arXiv preprint arXiv:1909.03683*, 2019.

[8] Bin Dai and David Wipf. Diagnosing and enhancing vae models. *arXiv preprint arXiv:1903.05789*, 2019.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[10] Aaron Fisher, Cynthia Rudin, and Francesca Dominici. All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously. *Journal of Machine Learning Research*, 20(177):1–81, 2019.

[11] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pages 3315–3323, 2016.

[12] Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. Women also snowboard: Overcoming bias in captioning models. In *European Conference on Computer Vision*, pages 793–811. Springer, 2018.

[13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.

[14] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. *Iclr*, 2(5):6, 2017.

[15] Daniel Im Jiwoong Im, Sungjin Ahn, Roland Memisevic, and Yoshua Bengio. Denoising criterion for variational auto-encoding framework. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.

[16] Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. Learning not to learn: Training deep neural networks with biased data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9012–9020, 2019.

[17] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018.

[18] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[19] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in neural information processing systems*, pages 4743–4751, 2016.

[20] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. *http://yann.lecun.com/exdb/mnist/*, 2010.

[21] Klas Leino, Emily Black, Matt Fredrikson, Shayak Sen, and Anupam Datta. Feature-wise bias amplification. *arXiv preprint arXiv:1812.08999*, 2018.

[22] Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9572–9581, 2019.

[23] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation bias. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 513–528, 2018.

[24] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015.

[25] Dominik Lorenz, Leonard Bereska, Timo Milbich, and Bjorn Ommer. Unsupervised part-based disentangling of object shape and appearance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10955–10964, 2019.

[26] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.

[27] Emile Mathieu, Tom Rainforth, N Siddharth, and Yee Whye Teh. Disentangling disentanglement in variational autoencoders. *arXiv preprint arXiv:1812.02833*, 2018.

[28] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016.

[29] Marie-Morgane Paumard, David Picard, and Hedi Tabia. Jigsaw puzzle solving using local feature co-occurrences in deep neural networks. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 1018–1022. IEEE, 2018.

[30] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. In *Advances in Neural Information Processing Systems*, pages 14837–14847, 2019.

[31] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in neural information processing systems*, pages 2234–2242, 2016.

[32] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.

[33] Rodrigo Santa Cruz, Basura Fernando, Anoop Cherian, and Stephen Gould. Deeppermnet: Visual permutation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3949–3957, 2017.

[34] Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. Ladder variational autoencoders. In *Advances in neural information processing systems*, pages 3738–3746, 2016.

[35] Tianbao Song, Jingbo Sun, Bo Chen, Weiming Peng, and Jihua Song. Latent space expanded variational autoencoder for sentence generation. *IEEE Access*, 7:144618–144627, 2019.

[36] Harini Suresh and John V Guttag. A framework for understanding unintended consequences of machine learning. *arXiv preprint arXiv:1901.10002*, 2019.

[37] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein autoencoders. *International Conference on Learning Representations*, 2018.

[38] Jakub M Tomczak and Max Welling. Vae with a vampprior. *arXiv preprint arXiv:1705.07120*, 2017.

[39] Serena Yeung, Anitha Kannan, Yann Dauphin, and Li Fei-Fei. Tackling over-pruning in variational autoencoders. *arXiv preprint arXiv:1706.03643*, 2017.

[40] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

[41] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. *arXiv preprint arXiv:1707.09457*, 2017.

[42] Qingyu Zhao, Nicolas Honnorat, Ehsan Adeli, Adolf Pfefferbaum, Edith V Sullivan, and Kilian M Pohl. Variational autoencoder with truncated mixture of gaussians for functional connectivity analysis. In *International Conference on Information Processing in Medical Imaging*, pages 867–879. Springer, 2019.