

IMPROVING THE UNSUPERVISED DISENTANGLED REPRESENTATION LEARNING WITH VAE ENSEMBLE

Anonymous authors

Paper under double-blind review

ABSTRACT

Variational Autoencoder (VAE) based frameworks have achieved the state-of-the-art performance on the unsupervised disentangled representation learning. A recent theoretical analysis shows that such success is mainly due to the VAE implementation choices that encourage a PCA-like behavior locally on data samples. Despite this implied model identifiability, the VAE based disentanglement frameworks still face the trade-off between the local orthogonality and data reconstruction. As a result, models with the same architecture and hyperparameter setting can sometimes learn entangled representations. To address this challenge, we propose a simple yet effective VAE ensemble framework consisting of multiple VAEs. It is based on the assumption that entangled representations are unique in their own ways, and the disentangled representations are “alike” (similar up to a signed permutation transformation). In the proposed VAE ensemble, each model not only maintains its original objective, but also encodes to and decodes from other models through pair-wise linear transformations between the latent representations. We show both theoretically and experimentally, the VAE ensemble objective encourages the linear transformations connecting the VAEs to be trivial transformations, aligning the latent representations of different models to be “alike”. We compare our approach with the state-of-the-art unsupervised disentangled representation learning approaches and show the improved performance.

1 INTRODUCTION

Disentangled representation learning aims to capture the semantically meaningful compositional representation of data (Higgins et al., 2018; Mathieu et al., 2018), and is shown to improve the efficiency and generalization of supervised learning (Locatello et al., 2019), reinforcement learning (Watters et al., 2019), and reasoning tasks (van Steenkiste et al., 2019). The current state-of-the-art unsupervised disentangled representation learning deploy the Variational Autoencoder (VAE) (Kingma & Welling, 2013; Rezende et al., 2014). The main challenge is to reduce the trade-off between learning a disentangled representation and reconstructing input data. Most of the recent works extend the original VAE objective with carefully designed augmented objective to address this trade-off (Higgins et al., 2017; Burgess et al., 2017; Kim & Mnih, 2018; Chen et al., 2018; Kumar et al., 2017). A recent study in (Locatello et al., 2018) compared these methods and showed that their performance is sensitive to initialization and hyperparameter setting of the augmented objective function.

Recently, Duan *et al.* (Duan et al., 2019) developed an unsupervised model selection method named Unsupervised Disentanglement Ranking (UDR) to address the challenge of hyperparameter search and model selection. UDR leverages the finding in (Rolinek et al., 2019) that the implementation choices of VAE encourage a local PCA-like behavior locally on data samples. As a result, disentangled representations by VAEs are “alike” as they are similar up to signed permutation transformations. On the contrary, the entangled representations by VAEs are “unique” as they are similar at least up to non-degenerate rotation matrices. UDR uses multiple models trained with different initializations and hyperparameter settings, and builds a similarity matrix measuring the pair-wise similarity between the latent variables from different models. A higher score is given to the model that can match its representations to many others models. The results show close match between UDR and commonly used supervised metrics, as well as the performance of downstream tasks using the latent representations.

45 Inspired by the findings from these studies, we propose a simple yet effective VAE ensemble frame-
 46 work to improve the disentangled representation by VAE. The proposed VAE ensemble consists of
 47 multiple VAEs. The latent variables in every pair of these models are connected through linear lay-
 48 ers to force the latent representations in the ensemble to be similar up to a linear transformation.
 49 We show that the VAE ensemble objective encourages these pair-wise linear transformations to con-
 50 verge to trivial transformations, making latent representations of different VAEs in the ensemble to
 51 be “alike”, thus disentangled. In this paper, we make the following contributions: (1) We introduce
 52 a simple yet effective VAE ensemble framework to improve the disentangled representation learning
 53 using the original VAE. (2) We show in theoretical analysis that the linear transformations connect-
 54 ing the latent representations of the individual models in the ensemble tend to converge to trivial
 55 transformations thus encourage disentangled representation, and verify this result with experiments.
 56 (3) We evaluate our approach using the original VAE model, and show the improved state-of-the-art
 57 performance across different datasets.

58 2 RELATED WORK

59 **Variational Autoencoder** is a deep directed probabilistic graphical model consisting of an encoder
 60 and a decoder (Kingma & Welling, 2013; Rezende et al., 2014). The encoder $q_\phi(z|x)$ maps the
 61 input data $x \in \mathbb{R}^n$ to a probabilistic distribution as the latent representation $z \in \mathbb{R}^d$, and the decoder
 62 $q_\theta(x|z)$ maps the latent representation to the data space noted as $q_\theta(x|z)$, where ϕ and θ represent
 63 model parameters. The VAE objective is to maximize the marginalized log-likelihood of data. Direct
 64 optimization of this objective is not tractable and it is approximated by the evidence lower bound
 65 (ELBO) as:

$$\mathcal{L}_{VAE} = \mathbb{E}_{q_\phi(z|x)}[\log q_\theta(x|z)] - \text{KL}(q_\phi(z|x) \parallel p(z)), \quad (1)$$

66 In practice, the first term is estimated by reconstruction error. The second term is the Kullback-
 67 Leibler divergence between the posterior $q_\phi(z|x)$ and the prior $p(z)$ commonly chosen as an
 68 isotropic unit Gaussian $p(z) \sim \mathcal{N}(0, \mathbf{I})$.

69 **Disentangled representation by VAE** has achieved the state-of-the-art performance (Higgins et al.,
 70 2017; Burgess et al., 2017; Kim & Mnih, 2018; Chen et al., 2018; Kumar et al., 2017), despite the
 71 fact that the VAE objective only models the marginal distribution of the data instead of the desired
 72 joint distribution over data and latent variables. The reason for this success is the implementation
 73 choices of the VAE framework (Rolinek et al., 2019). In practice, the latent variables in VAE often
 74 work in “polarized” modes. The “passive” mode is defined by $\mu_j^2(x) \ll 1$ and $\sigma_j^2(x) \approx 1$, while the
 75 “active” mode is defined by $\sigma_j^2(x) \ll 1$. The “passive” latent variables closely approximate the prior
 76 and have little effect on the decoder. The “active” latent variables, on the other hand, are closely
 77 related to both the per sample KL loss and the decoder output. The “polarized regime” enables a
 78 reformulated VAE objective showing that VAEs optimize a trade-off between data reconstruction
 79 and orthogonality of the linear approximation of decoder Jacobian locally around a data sample.
 80 This PCA-like behavior near data points encourages an identifiable disentangled latent space by
 81 VAE. Furthermore, it was suggested that finding an appropriate “polarized regime” is dependent
 82 on the initialization and the hyperparameter tuning of the state-of-the-art approaches. In this study,
 83 we show that the proposed VAE ensemble aligns the “polarized regime” of individual VAE models
 84 towards the disentangled representation.

85 **Model selection** In practice, we often observe neural networks achieve similar performance with
 86 different internal representations when trained with the same hyperparameters (Raghu et al., 2017;
 87 Wang et al., 2018; Morcos et al., 2018). For the unsupervised disentangled representation, as dis-
 88 cussed in (Locatello et al., 2018; Duan et al., 2019), we often observe high variance in the perfor-
 89 mance from the model trained with the same architecture and hyperparameter setting. This poses
 90 a challenge for choosing the model in practice. Duan et al. (2019) proposed Unsupervised Disen-
 91 tanglement Ranking (UDR) to address this challenge. The extensive empirical evaluations on UDR
 92 using both the supervised metric measurement and the performance of downstream tasks validates
 93 its effectiveness. They also confirm that disentangled representations are “alike” and entangled rep-
 94 resentations are unique in their own ways. The proposed VAE ensemble leverages this finding.

95 **Identifiable VAE** Built on the recent breakthroughs in nonlinear Independent Component Analy-
 96 sis (ICA) literature (Hyvarinen & Morioka, 2016; 2017; Hyvarinen et al., 2019), Khemakhem *et*
 97 *al.* show that the identification of the true joint distribution over observed and latent variables is

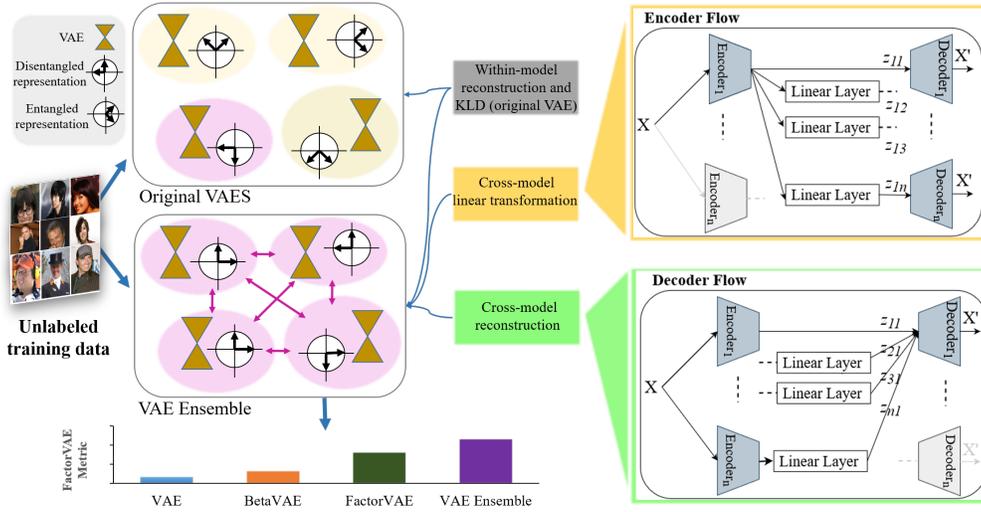


Figure 1: The proposed VAE ensemble consists of multiple original VAE models. The encoders of the VAEs in the ensemble generate input encoding that can be linearly transformed among each other. The decoders of the VAEs in the ensemble reconstruct the input data from both their corresponding encoder and the linearly transformed encodings from other encoders. The x and y axis of the circles on the left hand side of the plot represent two generative factors as an example. The aligned arrows with x and y axis show a model with disentangled representation and unaligned ones show a model with entangled representation.

possible up to very simple transformations (Khemakhem et al., 2019). They proposed identifiable VAE (iVAE) that requires a factorized prior distribution over the latent variables conditioned on an additionally observed variable, such as a class label or almost any other observation. We believe the proposed VAE ensemble is related to such framework where the latent representation from one VAE model can be regarded as the auxiliary observations for another.

Ensemble learning The idea of ensemble learning is to combine multiple learning models (potentially weak learners) to improve the task performance or robustness over a single model (Schapire, 1990). It achieves the improved performance by averaging the bias, reducing the variance thus preventing the over-fitting (Drucker et al., 1994; Breiman, 1996). Early works in neural networks have used the ensemble learning to achieve top performance in the related competition (Krizhevsky et al., 2012; Simonyan & Zisserman, 2014). In this work, we apply the ensemble learning to enforce the alignment among the latent representations of different models. This results in latent representations that are similar among each other in the ensemble up to a trivial transformation.

3 THE VAE ENSEMBLE FRAMEWORK

As illustrated in Figure 1, the proposed VAE ensemble consists of n original VAE models with the same architecture but different initializations. It also consists of $n \times (n - 1)$ linear layers connecting the latent representations of every two VAE models. Each model in the ensemble maintains its original VAE objective as Eq. 1. In addition, l_2 loss is used to force mapping between latent representations via pair-wise linear layers (cross-model linear transformation). The decoder of each VAE model generates the input reconstruction from not only their corresponding encoder (within-model reconstruction), but also the linearly transformed encodings from other encoders (cross-model reconstruction). Overall, the VAE ensemble is trained with the following objective:

$$\begin{aligned} \mathcal{L}(\theta, \phi) = & \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_{q_{\phi_{ij}}(z_{ij}|x)} [\log q_{\theta_j}(x|z_{ij})] - \sum_{i=1}^n \text{KL}(q_{\phi_{ii}}(z_{ii}|x) \parallel p(z_{ii})) \\ & - \gamma \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}_{q_{\phi_{ij}}(z_{ij}|x)} \|z_{jj} - z_{ij}\|^2, \end{aligned} \quad (2)$$

120 where n is the number of models in the ensemble; $\phi := (\phi_{ij})$ is the encoders parameters where ϕ_{ij}
 121 represents the encoder of VAE_{*i*} and its associated linear layer mapping the latent representation from
 122 VAE_{*i*} to VAE_{*j*} (notice that ϕ_{ii} represents the encoder parameters of VAE_{*i*} only and its associated
 123 linear transformation can be regarded as an identity transformation); $\theta := (\theta_i)$ represents the de-
 124 coder parameters of VAE_{*i*}; and z_{jj} represents the latent representation of VAE_{*j*} while z_{ij} represents
 125 the linearly transformed latent representation from VAE_{*i*} to VAE_{*j*}. γ is a hyperparameter to balance
 126 the effect of the estimation error between the latent representations. $p(z_{ii}) \sim \mathcal{N}(0, I)$ is assume to
 127 be the prior as defined in the original VAE objective.

128 Comparing to the original VAE objective in Eq. 1, the objective of each individual VAE model in
 129 the ensemble, \mathcal{L}'_{VAE} , can be written as:

$$\mathcal{L}'_{VAE}(\theta, \phi) = \mathcal{L}_{VAE}(\theta, \phi) + \underbrace{\sum_{j=1}^{n-1} \left\{ \mathbb{E}_{q_{\phi_j}(z_j|x)}[\log q_{\theta}(x|z_j)] - \gamma \mathbb{E}_{q_{\phi_j}(z_j|x)} \|z_j - z\|^2 \right\}}_{\text{Ensemble Regularization}}, \quad (3)$$

130 where n is the number of models in the ensemble, ϕ_j stands for the parameters of the encoder and
 131 its linear transformation layers from other VAEs, and z_j represents the linear transformed latent
 132 representation of the encoding from other encoders.

133 In this form, the VAE ensemble regularizes each VAE model with additional terms on the encoder
 134 as $\gamma \mathbb{E}_{q_{\phi_j}(z_j|x)} \|z_j - z\|^2$, and on the decoder as $\sum_{j=1}^n \mathbb{E}_{q_{\phi_j}(z_j|x)} [\log q_{\theta}(x|z_j)]$. These regulariza-
 135 tions directly constrain the latent representations among different VAE models in the ensemble to
 136 be similar. In particular, for a given input data, $\|z_j - z\|^2$ encourages the encoders to generate
 137 similar encodings up to the linear transformations; and $\sum_{j=1}^n \mathbb{E}_{q_{\phi_j}(z_j|x)} [\log q_{\theta}(x|z_j)]$ emphasizes
 138 the similar effect on the data reconstruction from the latent variables such that the decoders can
 139 reconstruct the input data with both the original encoding z and the linearly transformed encoding
 140 z_j . As we shall discuss in the next section, together these regularizations encourage similar latent
 141 representation up a trivial transformation by different models in the ensemble.

142 We also introduce the hyperparameter γ to balance the trade-off between these two regularizations:
 143 higher value forces closer mapping between the encoders and reduce the cross-model reconstruction
 144 error of the decoders; lower value relaxes the mapping between the encoders and increases the cross-
 145 model reconstruction error of the decoders. As we show in Section 5, both components are important
 146 and balancing the trade-off between them is important as the ensemble size increases.

147 **Computational complexity** It is a common practice to train a number of seeds per hyperparameter
 148 setting for the current state-of-the-art unsupervised disentanglement VAE models (Locatello et al.,
 149 2018; Duan et al., 2019). Comparing to training n original VAEs, the proposed VAE ensemble
 150 requires additional $n \times (n - 1)$ linear layers. While this addition does not increase the size of the
 151 model much, the estimation of the linear transformations loss and the cross-model reconstruction
 152 losses grow with $n \times (n - 1)$, which may be computationally expensive especially when n is large.
 153 That being said, the results in Section 5 show that the VAE ensemble achieves more stable results
 154 comparing to the current state-of-the-art models. Also, its computation is highly parallelisable.

155 4 THEORETICAL JUSTIFICATION

156 In this section, we present the theoretical analysis on why the proposed VAE ensemble can improve
 157 the disentangled representation. We start with analysing the l_2 objective in Eq. 2 of the pair-wise
 158 liner transformations in the VAE ensemble, and show that: (1) the pair-wise linear transformations
 159 encourage similar “polarized” regime (see Sec. 2) among the VAEs in the ensemble; (2) the linear
 160 transformations are close to the orthogonal transformations. Based on these two properties, we then
 161 discuss how the cross-model reconstructions by the VAE ensemble encourage learning a disentan-
 162 gled representation over its entangled counterpart.

163 4.1 THE EFFECT OF LINEAR TRANSFORMATION BETWEEN LATENT REPRESENTATIONS

164 Let VAE_{*i*} and VAE_{*j*} be two different VAE models in the ensemble, and M_{ji} be the linear trans-
 165 formation that maps the latent representation of a given input x by VAE_{*j*} to the one by VAE_{*i*}, as

$\mathbf{z}_j(x) \sim \mathcal{N}(\boldsymbol{\mu}_j(x), \text{diag}(\boldsymbol{\sigma}_j(x)^2))$ to $\mathbf{z}_i(x) \sim \mathcal{N}(\boldsymbol{\mu}_i(x), \text{diag}(\boldsymbol{\sigma}_i(x)^2))$. In the following we remove the input notation from the VAE latent representations for simplicity (i.e. $\mathbf{z}_j(x)$ is simplified as \mathbf{z}_j), while keeping in mind that the analysis is based on the local latent representation of a given input x .

For VAE_j , the l_2 term of the VAE ensemble loss in Eq. 2 aims to find M_{ji} and \mathbf{z}_j that minimize $\mathbb{E} \|\mathbf{z}_i - M_{ji}\mathbf{z}_j\|^2$, where the expectation is over the stochasticity of VAE_j . We can write \mathbf{z}_i and \mathbf{z}_j as $\mathbf{z}_i = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_i$ and $\mathbf{z}_j = \boldsymbol{\mu}_j + \boldsymbol{\epsilon}_j$, where $\boldsymbol{\epsilon}_i \sim \mathcal{N}(0, \text{diag}(\boldsymbol{\sigma}_i^2))$ and $\boldsymbol{\epsilon}_j \sim \mathcal{N}(0, \text{diag}(\boldsymbol{\sigma}_j^2))$. Hence using bias-variance decomposition, the l_2 term can be written as:

$$\begin{aligned} & \min_{M_{ji}, \mathbf{z}_j \sim \mathcal{N}(\boldsymbol{\mu}_j, \text{diag}(\boldsymbol{\sigma}_j^2))} \mathbb{E} \|\mathbf{z}_i - M_{ji}\mathbf{z}_j\|^2 \\ &= \min_{M_{ji}, \mathbf{z}_j \sim \mathcal{N}(\boldsymbol{\mu}_j, \text{diag}(\boldsymbol{\sigma}_j^2))} \left\{ \|\boldsymbol{\mu}_i - M_{ji}\boldsymbol{\mu}_j\|^2 + \mathbb{E}_{\mathbf{z}_j} \|M_{ji}\boldsymbol{\mu}_j - M_{ji}\mathbf{z}_j\|^2 + \mathbb{E}_{\mathbf{z}_i} \|\boldsymbol{\mu}_i - \mathbf{z}_i\|^2 \right\} \\ &= \min_{M_{ji}, \mathbf{z}_j \sim \mathcal{N}(\boldsymbol{\mu}_j, \text{diag}(\boldsymbol{\sigma}_j^2))} \left\{ \|\boldsymbol{\mu}_i - M_{ji}\boldsymbol{\mu}_j\|^2 + \mathbb{E}_{\mathbf{z}_j} \|M_{ji}\boldsymbol{\mu}_j - M_{ji}\mathbf{z}_j\|^2 \right\} + C_1 \\ &= \min_{M_{ji}, \mathbf{z}_j \sim \mathcal{N}(\boldsymbol{\mu}_j, \text{diag}(\boldsymbol{\sigma}_j^2))} \left\{ \|\boldsymbol{\mu}_i - M_{ji}\boldsymbol{\mu}_j\|^2 + \mathbb{E}_{\boldsymbol{\epsilon}_j} \|M_{ji}\boldsymbol{\epsilon}_j\|^2 \right\} + C_1 \end{aligned} \quad (4)$$

where the constant C_1 arises from the fact $\mathbb{E}_{\mathbf{z}_i} \|\boldsymbol{\mu}_i - \mathbf{z}_i\|^2$ does not depend on M_{ji} and \mathbf{z}_j . Eq. 4 consists of a deterministic component of $\|\boldsymbol{\mu}_i - M_{ji}\boldsymbol{\mu}_j\|^2$ and a stochastic component of $\mathbb{E}_{\boldsymbol{\epsilon}_j} \|M_{ji}\boldsymbol{\epsilon}_j\|^2$.

The deterministic component can be minimized by adjusting the parameters in VAE_j such that its mean encoding $\boldsymbol{\mu}_j$ is optimized for any given M_{ji} . This simplifies our analysis to focus on the stochastic component. Between M_{ji} and $\boldsymbol{\epsilon}_j$ in this stochastic component, we separately optimize one while having the other fixed.

We start with fixed M_{ji} and optimizing for $\boldsymbol{\epsilon}_j \sim \mathcal{N}(0, \text{diag}(\boldsymbol{\sigma}_j^2))$. Notice that $\boldsymbol{\sigma}_j^2$ is associated with VAE_j objective. In (Rolinek et al., 2019), the VAE objective is reformulated into the deterministic reconstruction, the stochastic reconstruction and the KL loss. The last two components define the stochastic loss of VAE. It is formulated as:

$$\min_{V, \boldsymbol{\sigma}_j^2} \sum_X \mathbb{E}_{\boldsymbol{\epsilon}_j \sim \mathcal{N}(0, \text{diag}(\boldsymbol{\sigma}_j^2))} \|D\boldsymbol{\epsilon}_j\|^2 \quad \text{s.t.} \quad \sum_X L_{KL} = C_1, \quad (5)$$

where X represents the dataset, D represents the local linear approximation of the Jacobian of the decoder with singular value decomposition as $D = U\Sigma V^T$. Furthermore, the KL loss $L_{KL} = \frac{1}{2} \sum_{k=1}^d (\mu_{jk}^2 + \sigma_{jk}^2 - \log \sigma_{jk}^2 - 1)$ can be simplified as $L_{\approx KL} = \frac{1}{2} \sum_{k \in \text{“active”}} (\mu_{jk}^2 - \log \sigma_{jk}^2 - 1)$ based on the “polarized” regime of VAE. (Rolinek et al., 2019) shows that $\boldsymbol{\sigma}_j^2$ act as the precision control allowed for each latent variable where more influential ones receive more precision. Combining the stochastic loss of the linear transformation in Eq. 4 and the stochastic loss of the original VAE in Eq. 5, the overall stochastic loss on $\boldsymbol{\sigma}_j^2$ can be formulated as:

$$\min_{\boldsymbol{\sigma}_j^2} \mathbb{E}_{\boldsymbol{\epsilon}_j \sim \mathcal{N}(0, \text{diag}(\boldsymbol{\sigma}_j^2))} [\|M_{ji}\boldsymbol{\epsilon}_j\|^2 + \|D\boldsymbol{\epsilon}_j\|^2] \quad \text{s.t.} \quad \sum_k -\log \sigma_{jk}^2 = C_2, \quad (6)$$

where σ_{jk}^2 is the k th element of $\boldsymbol{\sigma}_j^2$. Here we further simplify L_{KL} with the $L_{\approx KL}$ up to additive constants C_2 when $\boldsymbol{\mu}_j$ is fixed. In addition to the precision control of $\boldsymbol{\sigma}_j^2$ on VAE_j , this objective also aims to find an optimal distribution of $\boldsymbol{\sigma}_j^2$ that aligns the “polarized regime” among different VAEs. To see why, let c_k be the k th column of M_{ji} , we then have $\mathbb{E} \|M_{ji}\boldsymbol{\epsilon}_j\|^2 = \sum_k \|c_k\|^2 \sigma_{jk}^2$. The Arithmetic-Mean-Geometric-Mean (AM/GM) inequality suggests that $\sum_k \|c_k\|^2 \sigma_{jk}^2 \geq n \left(\prod_k \|c_k\|^2 \sigma_{jk}^2 \right)^{1/n} = n \left(\prod_k \|c_k\|^2 \right)^{1/n} \exp(-C)$, where the equality is achieved when $\|c_m\|^2 \sigma_{jm}^2 = \|c_n\|^2 \sigma_{jn}^2$ for any $m \neq n$. This suggests that latent variables with high $\|c_k\|^2$ mapping from \mathbf{z}_j to \mathbf{z}_i will have smaller variance. Hence, these latent variables in \mathbf{z}_j are encouraged to stay in the “active” mode. On contrary, the latent variables that do not share similar generative factors between \mathbf{z}_j and \mathbf{z}_i will be assigned larger variance, and being pushed towards the “passive” mode.

202 Now we fix the optimal distribution of σ_j^2 , and optimize for M_{ji} . Since $\epsilon_j \sim \mathcal{N}(0, \text{diag}(\sigma_j^2))$, this
 203 objective can be understood as optimally rotating the latent space of VAE_j such that the stochastic
 204 component in Eq. 4 is minimized. Specifically, we have the following objective:

$$\min_{M_{ji}} \|M_{ji}\epsilon_j\|^2 = \min_R \|M_{ji}R^T\epsilon_j\|^2 \quad (7)$$

205 where R is an orthogonal transformation. Let c'_k be the k th column of $M_{ji}R$. Similar as before, the
 206 AM/GM inequality suggests $\|M_{ji}R^T\epsilon_j\|^2 = \sum_k \|c'_k\|^2 \sigma_{jk}^2 \geq \prod_k \|c'_k\|^2 \exp(-C_3)$. Hadamard’s
 207 inequality suggests that $\prod_k \|c'_k\|^2 \geq |\det(M_{ji}R)|$, and the equality is satisfied when c'_k s are pair-
 208 wise orthogonal. This can be understood from the geometric perspective where $\prod_k \|c'_k\|^2$ gives an
 209 upper bound on $\text{Volume}(\{M_{ji}R^T x : x \in [0, 1]^d\})$. As a result, the optimization of M_{ji} will lead to
 210 an orthogonal transformation. Together the optimization of Eq. 6 and Eq. 7 encourages the align-
 211 ment of the “polarized regime” among different models under orthogonal linear transformations.
 212 They force different models in the ensemble to capture the same mixture of the generative factors.
 213 In the next section, we discuss the effect of the cross-model reconstruction in the VAE ensemble that
 214 encourages the disentangled representation over the entangled ones.

215 4.2 THE EFFECT OF CROSS-MODEL RECONSTRUCTION

216 In an entangled representation, each latent variable captures a mixture of generative factors in its
 217 unique way. Since different generative factors typically have different effects on data variations
 218 (Duan et al., 2019), the orthogonal transformation from one entangled representation \mathbf{z}_j to another
 219 one \mathbf{z}_i introduces different encoding variance. Some of the transformed latent variables in $M_{ji}\mathbf{z}_j$
 220 carry larger variance comparing to the corresponding ones in \mathbf{z}_i . This discrepancy leads to larger
 221 cross-model reconstruction of VAE_i than the within model reconstruction. This error forces both
 222 VAE_i and VAE_j to adjust their representations until the effect on the data reconstruction by indi-
 223 vidual latent variables matches between $M_{ji}\mathbf{z}_j$ and \mathbf{z}_i . The process applies to all models in the
 224 ensemble and eventually leads to a one-to-one mapping of latent variables between different mod-
 225 els, where M_{ji} becomes a trivial transformation (signed permutation matrix). In particular, if one of
 226 the models in the VAE ensemble learns a disentangled representation, other models in the ensemble
 227 will converge to it. This is because the orthogonal transformation from an entangled representation
 228 to a disentangled representation introduces larger cross-model encoding variance due to the mixture
 229 of different generative factors in the former, thus a larger cross-model reconstruction by the disen-
 230 tangled model. On contrary, the orthogonal transformation from a disentangled representation to an
 231 entangled representation would not introduce larger cross-model encoding variance than the within
 232 model encoding, thus similar cross-model reconstruction as within model reconstruction by the en-
 233 tangled model. Such a gap encourages the entangled representations to align with the disentangled
 234 representation. We illustrate the geometric interpretation of such a case in Appendix C.

235 From these discussions, we conclude that the VAE ensemble encourages different individual mod-
 236 els to capture similar generative factors, thus learn representations that are “alike” up to a trivial
 237 transformation. In the next section, we verify these analytic results with experiments.

238 5 EXPERIMENTS

239 Our experiments are designed to confirm the discussions in the previous sections. Particularly we
 240 ask the following questions: **Q1:** Do the linear transformations in the ensemble converge to trivial
 241 transformation? **Q2:** Do the VAEs in the ensemble work in similar “polarized” regime? **Q3:** Does
 242 VAE ensemble improve the unsupervised disentangled representation learning, and what is the ef-
 243 fect of ensemble size? **Q4:** What are the effects of the cross-model reconstruction loss, the linear
 244 transformation loss and the hyperparameter γ in the VAE ensemble objective?

245 We analyze the inner working of the proposed VAE ensemble using the benchmark *dSprite* dataset
 246 (Matthey et al., 2017) with fully known generative processes, and the real-world *CelebA* dataset (Liu
 247 et al., 2015) with unknown generative process. Furthermore, for *dSprites* dataset, we compare our
 248 results with the original VAE model and the state-of-the-art disentanglement VAE models including
 249 β -VAE (Higgins et al., 2017), FactorVAE (Kim & Mnih, 2018), TC-VAE (Chen et al., 2018) and
 250 DIP-VAE (Kumar et al., 2017). We use two widely used supervised metrics including FactorVAE
 251 metric (Kim & Mnih, 2018) and DCI Disentanglement scores (Eastwood & Williams, 2018) as the

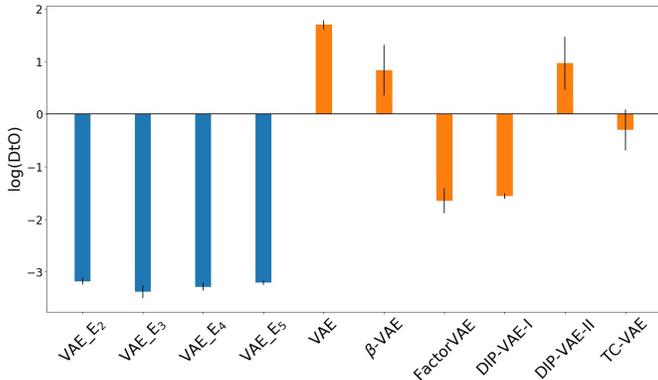


Figure 2: Comparing the DtO of linear transformations in the VAE ensemble ($\gamma=10$) with the one between well-trained individual VAEs, as well as the well-trained individual state-of-the-art VAE models. The latent dimension for all models is set to 10 and evaluated on the *dSprite* dataset.

quantitative measurements. They are shown to correlate with other common supervised metrics (Locatello et al., 2018). For example, FactorVAE metric and β -VAE metric (Higgins et al., 2017) capture similar notions, while DCI Disentanglement and Mutual Information Gap (MIG) (Chen et al., 2018) capture similar notions. In addition, DCI Disentanglement is closely related to the unsupervised model selection method UDR (Duan et al., 2019). For *CelebA* dataset, we show the latent traversal visualization as a qualitative measurement in Appendix E. We provide the details of the experiments in Appendix D.

Q1: We use the *Distance to Orthogonality* (DtO) (Rolinek et al., 2019) to check if the linear transformations in the ensemble converge to a signed permutation matrix during training. DtO is the Frobenius norm of the difference between a matrix M and its closest signed permutation matrix $P(M)$. It is solved with mixed-integer linear programming (MILP) formulation. The details on DtO can be found in Appendix B. In Figure 2, we show the DtO estimation of the linear transformations in the VAE ensemble of different ensemble size for the *dSprite* dataset. We show the mean and standard deviation of DtO across all linear transformations over 10 different runs. Furthermore, we compare these results with a baseline measurement where DtO is estimated for the linear transformations between the mean latent representations of well-trained individual models. Specifically, we use ten well-trained individual models and report the mean and standard deviation of the DtO estimations. As seen in the figure, the VAE ensemble models with different ensemble size all approach to trivial transformations between the individual models, while other VAE models do not have such property. In Fig. 6, we show that during training, the VAE ensemble remains maintains low DtO while the original VAEs do not have such property. A similar result for models trained on the *CelebA* dataset with different latent dimensions is shown in Fig. 7. Further discussion on these results are provided in Appendix E.

Q2: To check if the models in the VAE ensemble work in similar “polarized” regime, we estimate the relative error between L_{KL} and $L_{\approx KL}$ as $\Delta = \frac{L_{KL} - L_{\approx KL}}{L_{KL}}$ for each latent variable. Smaller Δ indicates closer matching between L_{KL} and $L_{\approx KL}$ of a latent variable, thus more “active”. Figure 3(a) and 3(b) show $\log(\Delta)$ of the 10 latent variables of individual models in VAE_E2 and VAE_E3 with different γ settings trained on the *dSprite* dataset. The results show that individual models in the ensemble do work in similar “polarized regime”. In Figure 3(a), we also compare the VAE ensemble with the β -VAE where $\beta = 4$. This setting was found previously to be the optimal setting for the *dSprite* data for β -VAE (Higgins et al., 2017). We see that the VAE ensemble encourages more “active” latent variables than β -VAE. When we compare Fig. 3(a) and 3(b), we see that as the ensemble size increases, individual models are forced to have more “active” latent variables by decomposing the generative factors. This can be observed in the latent traversals shown in Appendix E. The *dSprites* dataset contains five ground truth generative factors. The VAE_E2 models can have up to eight “active” latent variables depending on input, and these representations capture a decomposition of the ground truth generative factors.

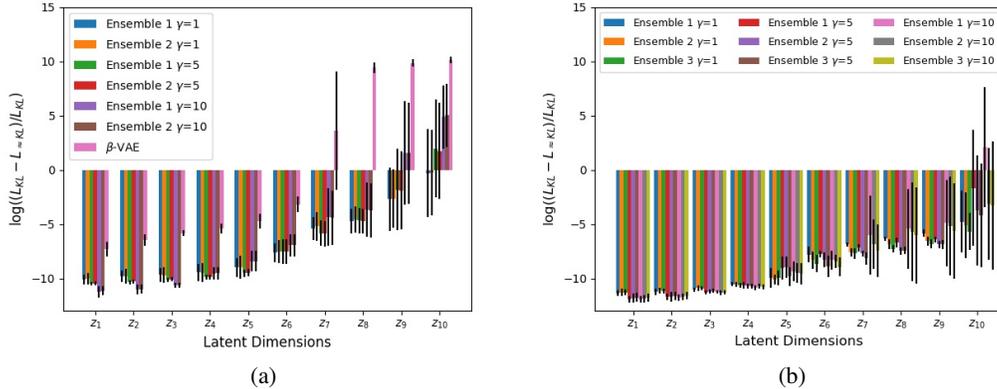


Figure 3: The “polarized regime” comparison between models in the VAE ensemble. The latent dimension is set to 10 and the results are over 10 runs of training on the *dSprite* dataset. (a) The “polarized regime” comparison by $\log(\frac{L_{KL} - L_{\approx KL}}{L_{KL}})$ of each latent variable of the two models in VAE_E₂ as well as β -VAE model. (b) Similar to (a) but for the three models in VAE_E₃.

FactorVAE Metric					
Individual Model		VAE Ensemble (VAE_E)			
VAE	0.635±0.083		$\gamma=1$	$\gamma=5$	$\gamma=10$
β -VAE ($\beta=4$)	0.665±0.089	VAE_E ₂	0.711±0.106	0.736±0.085	0.741±0.086
FactorVAE ($\gamma=40$)	0.764±0.075	VAE_E ₃	0.794±0.030	0.792±0.075	0.821±0.066
DIP-VAE-I ($\lambda_{od}=5$)	0.638±0.108	VAE_E ₄	0.833±0.037	0.790±0.038	0.800±0.078
DIP-VAE-II ($\lambda_{od}=5$)	0.676±0.122	VAE_E ₅	0.828±0.016	0.786±0.051	0.739±0.085
TC-VAE ($\beta=4$)	0.808±0.079				

DCI-Disentanglement Metric					
Individual Model		VAE Ensemble (VAE_E)			
VAE	0.143±0.033		$\gamma=1$	$\gamma=5$	$\gamma=10$
β -VAE ($\beta=4$)	0.198±0.076	VAE_E ₂	0.176±0.043	0.243±0.029	0.201±0.037
FactorVAE ($\gamma=40$)	0.253±0.072	VAE_E ₃	0.214±0.064	0.236±0.051	0.311±0.060
DIP-VAE-I ($\lambda_{od}=5$)	0.049±0.017	VAE_E ₄	0.240±0.059	0.223±0.045	0.251±0.038
DIP-VAE-II ($\lambda_{od}=5$)	0.106±0.032	VAE_E ₅	0.242±0.032	0.244±0.039	0.196±0.050
TC-VAE ($\beta=4$)	0.303±0.052				

Table 1: Comparison between the proposed VAE ensemble, the original VAE, and the current state-of-the-art disentangled VAE models. We report the mean and standard deviation of the FactorVAE metric and DCI Disentanglement scores over 10 runs trained on the *dSprite* data.

289 **Q3:** In Table 1, we compare the disentangled representation performance between the proposed
290 VAE ensemble and the state-of-the-art models. For the VAE ensemble, we report the performance
291 of the first model in the ensemble. We also report the results for the VAE ensemble with different
292 ensemble size and γ values. As shown in the table, the VAE ensemble significantly improves the
293 performance over the original VAE model. In many settings, the VAE ensemble achieves similar
294 or better performance over the state-of-the-art models. In Table 2, we evaluate the consistency
295 among the models in the ensemble by reporting the standard deviation of the evaluation metrics
296 using different models in the same ensemble. The small values confirm that different models in
297 the ensemble learn similar latent representations. Furthermore, Table 1 shows the joint effect of
298 ensemble size and γ setting. When $\gamma = 1$, the performance of VAE ensemble increases as the
299 ensemble size increases, indicated by the higher mean and smaller variance of both the FactorVAE
300 and DCI Disentanglement metrics. This behavior is consistent with the characteristic of ensemble
301 learning where the increase in performance becomes smaller as the size of ensemble increases.
302 However, as γ increases, having larger ensemble size can reduce the performance. We believe this

	γ	VAE_E ₂	VAE_E ₃	VAE_E ₄	VAE_E ₅
FactorVAE Metric	($\gamma=1$)	0.0019	0.0090	0.0048	0.0081
	($\gamma=5$)	0.0058	0.0064	0.0089	0.0163
	($\gamma=10$)	0.0060	0.0046	0.0147	0.0139
DCI-Disent Metric	($\gamma=1$)	0.0024	0.0026	0.0028	0.0036
	($\gamma=5$)	0.0037	0.0054	0.0049	0.0040
	($\gamma=10$)	0.0013	0.0024	0.0041	0.0042

Table 2: The comparison between individual models in the same ensemble. We report the average of the standard deviation of the metrics by individual models in the ensemble across 10 runs.

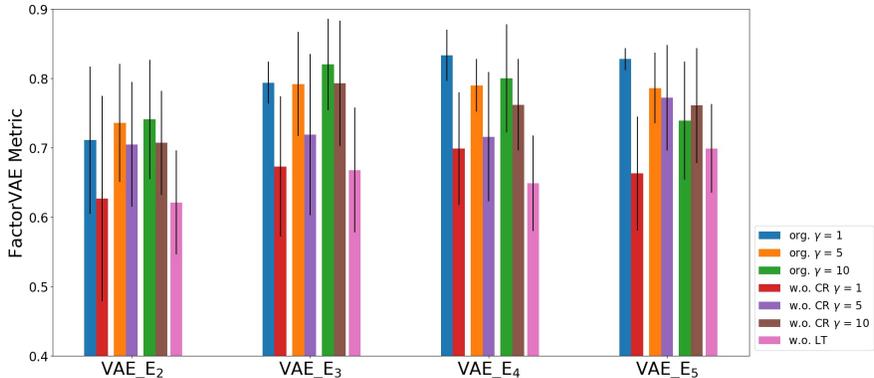


Figure 4: Ablation study to understand the effect of cross-model reconstruction and linear transformation in the VAE ensemble objective using the FactorVAE metric. (*w.o. CR* - without cross-model reconstruction loss; *w.o. LT* - without linear transformation loss; *org* - original VAE ensemble loss)

is due to the increased difficulty of balancing between the cross-model and within model objectives of VAE ensemble for larger ensembles. The reduced alignment of the latent representations among different models can also be seen in Table 2 where difference in the performance among individual models in the ensemble increases as ensemble size increases.

Q4: We conduct the ablation study to further understand the effect of the linear transformation loss and the cross-model reconstruction loss in the VAE ensemble objective. As shown in Fig. 4, removing either component leads to a lower FactorVAE metric for the VAE ensemble. Without the linear transformation loss, the performance of VAE ensemble decreases significantly across different ensemble sizes. Without the cross-model reconstruction loss, the performance of VAE ensemble also decreases but the gap becomes smaller as γ increases. This matches the discussion in Section 3 that higher γ forces closer mapping between the encoders and reduce the cross-model reconstruction error of the decoders. However, this also reduces the effect of cross-model reconstruction as discussed in Section 4.2. A similar result is also found for the DCI Disentanglement metric as shown in Fig. 8 in Appendix E. Overall, the results from the ablation study confirms the importance of both the linear transformation loss and the cross-model reconstruction loss in the VAE ensemble objective.

6 CONCLUSION

In this study, we propose a simple yet effective VAE ensemble framework consisting of multiple original VAEs to learn disentangled representation. The individual models in the ensemble are connected through linear layers that regularize both encoders and decoders to align the latent representations to be similar up to a linear transformation. We show in theory that the regularization by the VAE ensemble forces the linear transformations to be trivial transformations and show improved performance on the unsupervised disentangled representation learning. The theoretical discussion in Section 4 is based on the original VAE objective, and our experiments also focus on the ensemble with original VAE. We believe such framework can be extended to other disentangled VAE models, or even a mixture of different VAE models, as long as the regularization by the ensemble does not conflict with the augmented objective of these models.

329 REFERENCES

- 330 Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- 331 Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins,
332 and Alexander Lerchner. Understanding disentangling in β -VAE. (NIPS), 2017.
- 333 Tian Qi Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating Sources of Disentangle-
334 ment in Variational Autoencoders. (ICLR):1–13, 2018.
- 335 Harris Drucker, Corinna Cortes, Lawrence D Jackel, Yann LeCun, and Vladimir Vapnik. Boosting
336 and other ensemble methods. *Neural Computation*, 6(6):1289–1301, 1994.
- 337 Sunny Duan, Loic Matthey, Andre Saraiva, Nick Watters, Chris Burgess, Alexander Lerchner, and
338 Irina Higgins. Unsupervised model selection for variational disentangled representation learning.
339 In *International Conference on Learning Representations*, 2019.
- 340 Cian Eastwood and Christopher K. I. Williams. A framework for the quantitative evaluation of
341 disentangled representations. *ICLR*, 2018.
- 342 Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick,
343 Shakir Mohamed, and Alexander Lerchner. β -Vae: Learning Basic Visual Concepts With a Con-
344 strained Variational Framework. *ICLR 2017*, (July):1–13, 2017.
- 345 Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and
346 Alexander Lerchner. Towards a Definition of Disentangled Representations. *arXiv*, pp. 1–29,
347 2018.
- 348 Aapo Hyvarinen and Hiroshi Morioka. Unsupervised feature extraction by time-contrastive learning
349 and nonlinear ica. In *Advances in Neural Information Processing Systems*, pp. 3765–3773, 2016.
- 350 Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and
351 generalized contrastive learning. In *The 22nd International Conference on Artificial Intelligence
352 and Statistics*, pp. 859–868, 2019.
- 353 AJ Hyvarinen and Hiroshi Morioka. Nonlinear ica of temporally dependent stationary sources.
354 *Proceedings of Machine Learning Research*, 2017.
- 355 Ilyes Khemakhem, Diederik P Kingma, and Aapo Hyvärinen. Variational autoencoders and nonlin-
356 ear ica: A unifying framework. *arXiv preprint arXiv:1907.04809*, 2019.
- 357 Hyunjik Kim and Andriy Mnih. Disentangling by factorising. *Proceedings of Machine Learning
358 Research*, 80:2649–2658, 10–15 Jul 2018.
- 359 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114,
360 2013.
- 361 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
362 lutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105,
363 2012.
- 364 Abhishek Kumar, Prasanna Sattigeri, and Avinash Balakrishnan. Variational inference of disentan-
365 gled latent concepts from unlabeled observations. *CoRR*, abs/1711.00848, 2017.
- 366 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.
367 In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 368 Francesco Locatello, Stefan Bauer, Mario Lucic, Sylvain Gelly, Bernhard Schölkopf, and L G Dec.
369 Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representa-
370 tions. *arXiv*, pp. 1–33, 2018.
- 371 Francesco Locatello, Gabriele Abbati, Thomas Rainforth, Stefan Bauer, Bernhard Schölkopf, and
372 Olivier Bachem. On the fairness of disentangled representations. In *Advances in Neural Informa-
373 tion Processing Systems*, pp. 14584–14597, 2019.

- Emile Mathieu, Tom Rainforth, N Siddharth, and Yee Whye Teh. Disentangling disentanglement in variational autoencoders. *arXiv preprint arXiv:1812.02833*, 2018. 374
375
- Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017. 376
377
- Ari Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural networks with canonical correlation. In *Advances in Neural Information Processing Systems*, pp. 5727–5736, 2018. 378
379
380
- Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. In *Advances in Neural Information Processing Systems*, pp. 6076–6085, 2017. 381
382
383
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014. 384
385
- Michal Rolinek, Dominik Zietlow, and Georg Martius. Variational autoencoders pursue pca directions (by accident). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 12406–12415, 2019. 386
387
388
- Robert E Schapire. The strength of weak learnability. *Machine learning*, 5(2):197–227, 1990. 389
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 390
391
- Sjoerd van Steenkiste, Francesco Locatello, Jürgen Schmidhuber, and Olivier Bachem. Are disentangled representations helpful for abstract visual reasoning? In *Advances in Neural Information Processing Systems*, pp. 14222–14235, 2019. 392
393
394
- Liwei Wang, Lunjia Hu, Jiayuan Gu, Zhiqiang Hu, Yue Wu, Kun He, and John Hopcroft. Towards understanding learning representations: To what extent do different neural networks learn the same representation. In *Advances in Neural Information Processing Systems*, pp. 9584–9593, 2018. 395
396
397
398
- Nicholas Watters, Loic Matthey, Matko Bosnjak, Christopher P Burgess, and Alexander Lerchner. Cobra: Data-efficient model-based rl through unsupervised object discovery and curiosity-driven exploration. *arXiv preprint arXiv:1905.09275*, 2019. 399
400
401

402 A TECHNICAL LEMMAS

403 In this section, we give the lemmas used in the theoretical discussion in Section 4.

404 **Lemma 1.** (*Jensen’s inequality*) *If $g(x)$ is a convex transformation on x , then this convex transfor-*
 405 *mation of a mean $g[\mathbb{E}(x)]$ is less than or equal to the mean of the convex transformed value $\mathbb{E}[g(x)]$;*
 406 *it is a simple corollary that the opposite is true of concave transformations.*

407 **Lemma 2.** (*AM-GM inequality*) *As an extension of Jensen’s inequality, given a list of non-negative*
 408 *real numbers x_1, x_2, \dots, x_n , the arithmetic mean of this list $\frac{1}{n} \sum_{i=1}^n x_i$ is greater than or equal*
 409 *to the geometric mean of the same list $(\prod_{i=1}^n x_i)^{\frac{1}{n}}$; and further, the equality holds if and only if*
 410 *$x_1 = x_2 = \dots = x_n$.*

411 **Lemma 3.** (*Hadamard’s inequality*). *if M is the matrix having columns c_i , then $|\det(M)| \leq$*
 412 *$\prod_{i=1}^n \|c_i\|$; and the equality in Hadamard’s inequality is achieved if and only if the vectors are*
 413 *orthogonal.*

414 B DISTANCE TO ORTHOGONALITY (DTO)

415 In this section, we introduce the detail of *Distance to Orthogonality* (DtO) that is used in our exper-
 416 iment to check if the linear transformations in the VAE ensemble approach trivial transformations.
 417 This measurement is also used in (Rolinek et al., 2019) for a similar purpose. DtO is the Frobe-
 418 nius norm of the difference between a square matrix M and its closest signed permutation matrix
 419 $P(M)$. Finding $P(M)$ can be formulated as a mixed-integer linear programming (MILP) problem
 420 as following:

$$\begin{aligned} \min_P \quad & \sum_{i,j} |M_{i,j} - P(M)_{i,j}| \\ \text{s.t.} \quad & P(M)_{i,j} \in \{-1, 0, 1\}, \quad \forall(i, j) \\ & \sum_i |P_{i,j}| = 1, \quad \forall j \\ & \sum_j |P_{i,j}| = 1, \quad \forall i \end{aligned} \quad (8)$$

421 By introducing new variables $P_{i,j}^+, P_{i,j}^- \in \{0, 1\}$ and $D_{i,j} = |M_{i,j} - P(M)_{i,j}|$, we can reformulate
 422 the above optimization problem as:

$$\begin{aligned} \min_P \quad & \sum_{i,j} D_{i,j} \\ \text{s.t.} \quad & (P_{i,j}^+ - P_{i,j}^-) - M_{i,j} \leq D_{i,j}, \quad \forall(i, j) \\ & M_{i,j} - (P_{i,j}^+ - P_{i,j}^-) \leq D_{i,j}, \quad \forall(i, j) \\ & \sum_i (P_{i,j}^+ + P_{i,j}^-) = 1, \quad \forall j \\ & \sum_j (P_{i,j}^+ + P_{i,j}^-) = 1, \quad \forall i \end{aligned} \quad (9)$$

423 Using this optimization formulation, DoT of a given matrix $M \in \mathbb{R}^{n \times n}$ is defined as:

$$DoT = \frac{1}{n^2} \sum_{i,j} |M_{i,j} - P(M)_{i,j}| \quad (10)$$

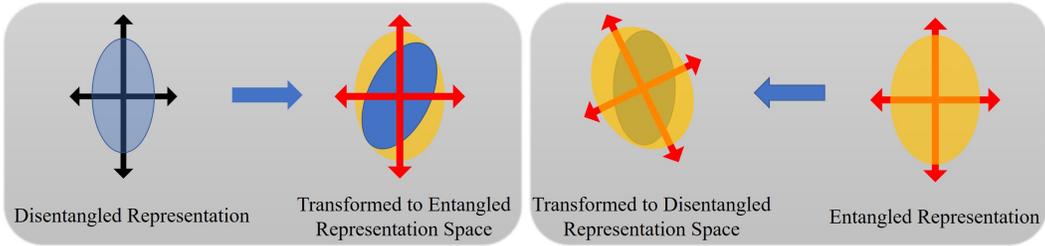


Figure 5: Geometric interpretation of the cross-model reconstruction between a disentangled representation space and an entangled representation space.

C GEOMETRIC INTERPRETATION OF THE EFFECT OF CROSS-MODEL RECONSTRUCTION 424
 RECONSTRUCTION 425

Given a disentangled and entangled latent representation space, Fig. 5 illustrates the effect of the cross-model reconstruction by VAE ensemble. The left part shows the orthogonal transformation from a disentangled representation to an entangled space, and the right part shows the transformation in the opposite direction. As shown in the figure, the orthogonal transformation from the disentangled representation to the entangled space does not introduce larger variance than the entangled representation. Hence, we can expect similar cross-model reconstruction and within model reconstruction. However, the transformation from the entangled representation to the disentangled space introduces larger variance (yellow shaded area over blue area on the right) than the disentangled representation. This leads to larger cross-model reconstruction by the disentangled model. 426
 427
 428
 429
 430
 431
 432
 433
 434

D MODEL ARCHITECTURE AND TRAINING DETAILS 435

We conducted our experiments, including training and evaluating the current state-of-the-art disentanglement models as well as evaluating the proposed VAE ensembles, using the `disentanglement_lib`¹ open-source library (Locatello et al., 2018). 436
 437
 438

Table 3 shows the encoder and the decoder architecture of the VAE model used in our experiments. This architecture is the same as the one used in the original β -VAE Higgins et al. (2017). 439
 440

Encoder	Decoder
Input 64×64 binary/RGB image	Input \mathbb{R}^d
4×4 conv, 32 ReLu, stride 2, pad 1	FC $d \times 256$, ReLu
4×4 conv, 32 ReLu, stride 2, pad 1	4×4 upconv, 64 ReLu, stride 1
4×4 conv, 64 ReLu, stride 2, pad 1	4×4 conv, 64 ReLu, stride 2, pad 1
4×4 conv, 64 ReLu, stride 2, pad 1	4×4 conv, 32 ReLu, stride 2, pad 1
4×4 conv, 256 ReLu, stride 1	4×4 conv, 32 ReLu, stride 2, pad 1
FC $256 \times (2 \times d)$	4×4 conv, nc , stride 2, pad 1

Table 3: Encoder and Decoder architecture, d : dimension of the latent representation; nc : number of input image channel (For $dSprites$ dataset $nc = 1$, for $CelebA$ dataset $nc = 3$).

Table 4 shows the hyperparameters setting used throughout the experiments. These parameters are fixed for all the experiments. 441
 442

E ADDITIONAL EXPERIMENTAL RESULTS 443

In this section, we present the additional results including the DtO and “polarized regime” analysis on the models trained on the *CelebA* dataset similar to the ones conducted on *dSprite* dataset in Section 5; the ablation results with DCI-Disentanglement metric and the DtO estimation; and the 444
 445
 446

¹https://github.com/google-research/disentanglement_lib

Parameter	value
Batch size	64
Latent dimension	10
Optimizer	Adam
Adam: beta1	0.9
Adam: beta2	0.999
Learning rate	1e-4

Table 4: Hyperparameters setting.

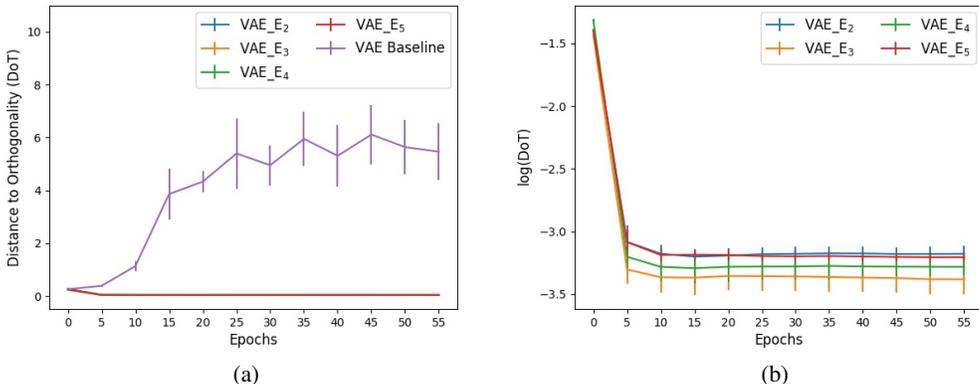


Figure 6: Characteristics of the linear transform between latent representations. The latent dimension is set to 10 and the results are over 10 runs of training on the *dSprite* dataset. (a) Comparing the DtO of linear transformations in the VAE ensemble ($\gamma=10$) and the one between original VAEs. (b) VAE ensemble ($\gamma=10$) with different ensemble size all achieve small DtO of the linear transformations between the models.

447 latent traversal of the trained model on both *dSprite* and *CelebA* dataset along with further discuss
 448 the effect of VAE ensemble on the latent representation.

449 E.1 CHARACTERIZATION OF THE LINEAR TRANSFORMATION IN VAE ENSEMBLE

450 In Figure 6, we show the DtO estimation of the linear transformations in the ensemble during training
 451 for the *dSprite* dataset. We report the mean and standard deviation of DtO across all linear
 452 transformations over 10 different runs. Furthermore, we compare these results with a VAE baseline
 453 where DtO is estimated for the linear transformations between original VAEs. Specifically, we train
 454 ten different VAEs separately and estimate the DtO of the pairwise linear transformation among
 455 these models during training. Similarly we report the mean and standard deviation of these DtO es-
 456 timations. As seen in the figure, the VAE ensemble models with different ensemble size all approach
 457 to trivial transformations between the individual models, while the original VAEs do not have such
 458 property. A similar result is also found in models trained for *CelebA* dataset. Similar to the results
 459 in Figure 6, we observe decreased DtO of the linear transformations in the VAE ensemble during
 460 training.

461 We also compare models trained with different latent dimension size. We observe decreased DtO
 462 as the latent dimension of the model increases in Figure 7. This is because, as discussed in the
 463 main paper, the VAE ensemble encourages more “active” latent variables. Models with higher latent
 464 dimension likely to learn a decomposition of generative factors. As a result, the alignment of the
 465 latent variable between different models are easier thus the linear transformations between the latent
 466 representations is closer to the trivial transformation. On the contrary when there are less latent
 467 variables in the model than the generative factors, some of the latent variables will capture more
 468 than a single generative factor. As a result, the one-to-one mapping between the latent variables of
 469 different models will not lead to a trivial transformation.

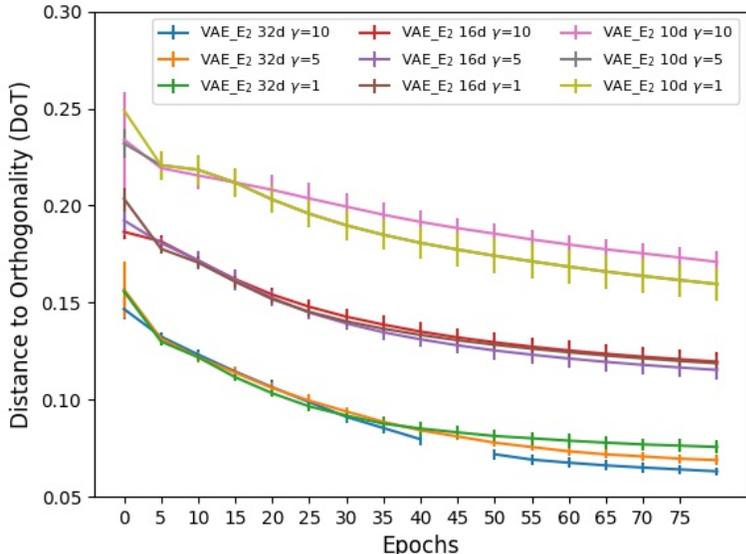


Figure 7: Distance to Orthogonality (DtO) measurement of the linear transforms between latent representations in VAE_E2 during training on the *CelebA* dataset. We also compare models with different latent dimensions of 10, 16 and 32 and the results are averaged over 5 runs. In the figure legend, we use “VAE_E_i nd $\gamma = g$ ” to represent VAE Ensemble (VAE.E) model with i individual VAE models, n latent dimensions and γ value equal to g .

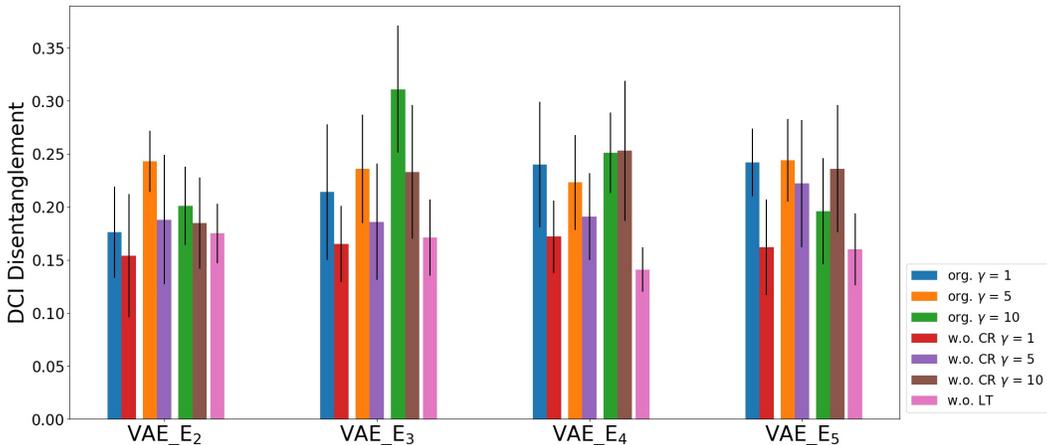
In Figure 9 and Figure 10, we show the “polarized regime” estimation for models in VAE_E2 and VAE_E3 trained for *CelebA* datasets, respectively. Similar to the results in Figure 3, individual models in the VAE ensemble tend to have similar “polarized regime”, and higher γ enforces the “polarized regime” by separating “passive” latent variables from the “active” ones. When we compare between VAE_E2 and VAE_E3, we observe increased “active” latent variables similar to the result on *dSprite* dataset in Section 5. More importantly, as discussed earlier, latent variables in a model with limited latent dimensions need to capture more than a single generative factor, especially for a complicated real-world dataset such as *CelebA*. This makes the linear transformation between the latent representations less trivial. As the latent dimension size grows, such constraint is relaxed and the linear transformations are closer to trivial.

These additional results confirm the conclusion in Section 5: (1) as the ensemble size increases, DtO increases due to the difficulty of aligning the latent representations among different models; (2) as the model latent dimension increases, DtO decreases due to the increased model capacity, and encourages the one-to-one mapping between latent variables in different models; (3) hyperparameter γ does not affect DtO significantly, but plays an important role on separating “active” and “passive” latent variables, especially when the latent dimension is large enough.

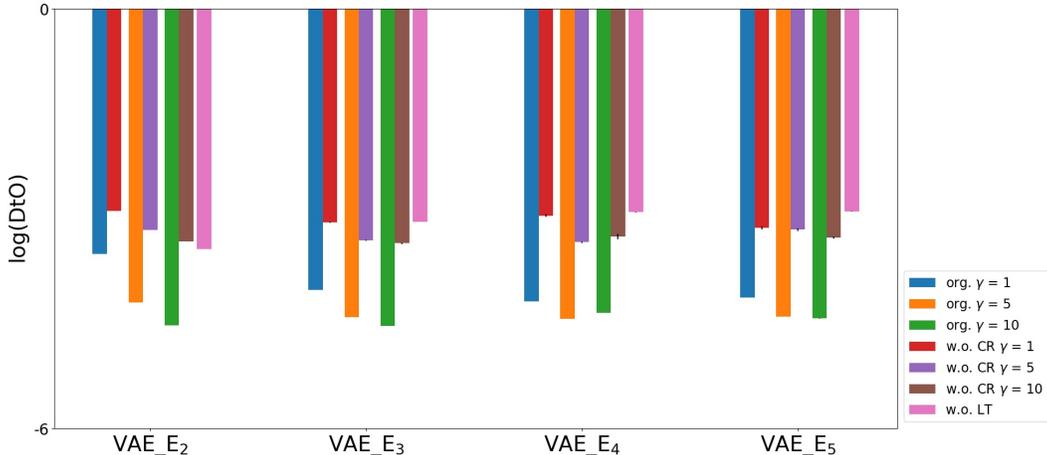
Furthermore, we believe the DtO measurement of the linear transformation in VAE ensemble could be a useful indicator for latent dimension size. As shown in Figure 7 and Figure 9, when the latent dimension is sufficient for a given dataset, the DtO of the linear transformation is small and some latent variables are pushed to “passive” mode.

E.2 ABLATION STUDY

Similar to the ablation result shown in Section 5, here we show the same ablation study using the DCI Disentanglement metric in Fig. 8(a) as well as the DtO measurement 8(b). Similar as the results of the FactorVAE metric in Fig. 4, removing either component leads to a lower DCI Disentanglement metric for the VAE ensemble. Without the linear transformation loss, the performance of VAE ensemble decreases significantly across different ensemble sizes. Fig. 8(b) shows that for VAE



(a) Comparison on the DCI-Disentanglement score



(b) Comparison on the Distance to Orthogonality (DtO)

Figure 8: Ablation study to understand the effect of cross-model reconstruction and linear transformation in the VAE ensemble objective using the DCI Disentanglement metric and DtO. (*w.o. CR* - without cross-model reconstruction loss; *w.o. LT* - without linear transformation loss; *org* - original VAE ensemble loss)

496 ensemble without the cross-model reconstruction, the linear transformations among models are close
 497 to a trivial transformation (signed permutation). This implies the orthonormal transformation of the
 498 linear transformations. This result further supports our intuitive justification in Appendix C that the
 499 cross-model objective encourages entangled models to align to disentangled models. Indeed, we see
 500 that adding the cross-model reconstruction can further reduce the DtO of the linear transformations
 501 among the models in the ensemble.

502 E.3 LATENT TRAVERSAL

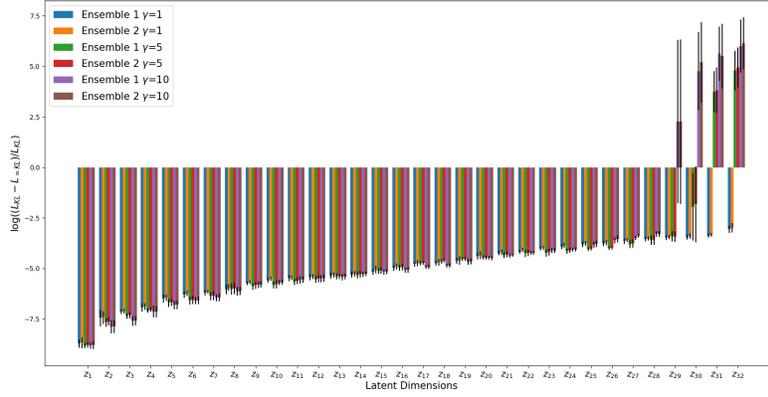
503 In this section we show the latent traversal of models trained on both *dSprites* and *CelebA* datasets.
 504 For a fixed input image, to extract the latent traversal we change the value of a single latent variable
 505 z_i in the corresponding encoding, and observe the generated output image to understand the effect
 506 of z_i . The range of the value are usually chosen to be from -3 to 3 due to the standard Gaussian
 507 prior.

508 In Figure 11, we show the latent traversal for both VAE_E2 and a single VAE model with 10 latent
 509 dimensions trained on *dSprites* dataset. Three images as shown in the last column of each block

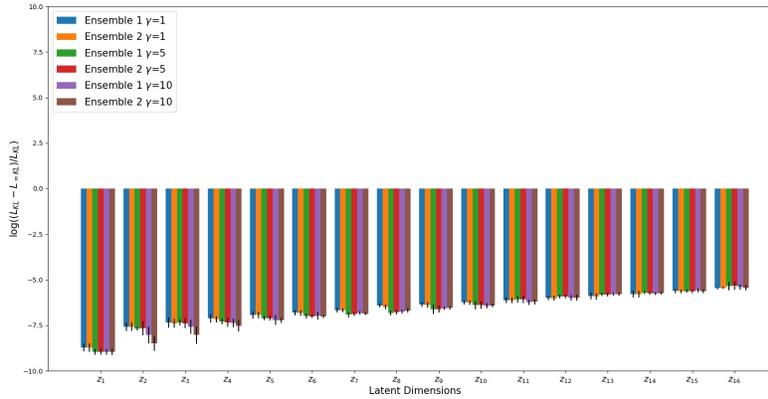
are used as input. Both models are able to capture certain generative factors of the data including position, shape, rotation and scale. In Section 5, we argue that the representation by VAE ensemble encourages more “active” latent variables, thus can capture a decomposition of the ground truth generative factors. Especially from the “polarized regime” estimation in Figure 3, we observe that some latent variables in the VAE ensemble are in-between “active” and “passive” modes. This suggests that the VAE ensemble model generates input-dependent factors based on the input complexity. In Figure 11, we observe this behavior highlighted with color boxes. The traversal on the second latent variable z_2 shows that an ellipse shape does not lead to an “active” latent variable. However, both heart and square shape lead to an “active” latent variable that changes the output. In contrast, the single VAE model does not have such behavior where the “active” modes are consistent across different input data.

In Figure 12 and Figure 13, we show the latent traversal for both VAE_E₂ and a single VAE model with 16 latent dimensions trained on *CelebA* dataset, respectively. In this real-world dataset, the generative factors are unknown. We observe different factors including background, azimuth, gender, hair style being captured by both models. Similar as before, the single VAE model maintains similar “active” mode for all latent variables where similar traversal patterns are observed for both input images. However, VAE_E₂ shows semantically consistent but input-dependent “active” mode. This is translated into different traversal effects and more realistic and sharper images by VAE_E₂, especially for the first input image that is less common in the dataset. We believe this is important towards a meaningful compositional latent representation learning.

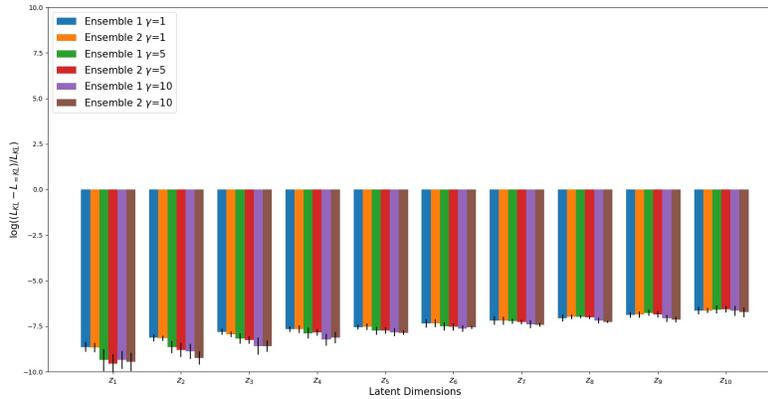
Overall, the latent traversal results in this section confirm the findings on the inner working of the VAE ensemble shown in the previous section as well as the discussion in Section 5.



(a) VAE_E2 with latent dimension of 32, CelebA data.

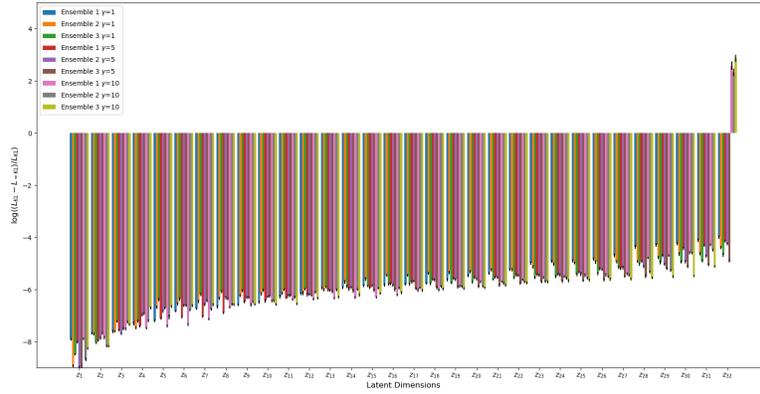


(b) VAE_E2 with latent dimension of 16, CelebA data.

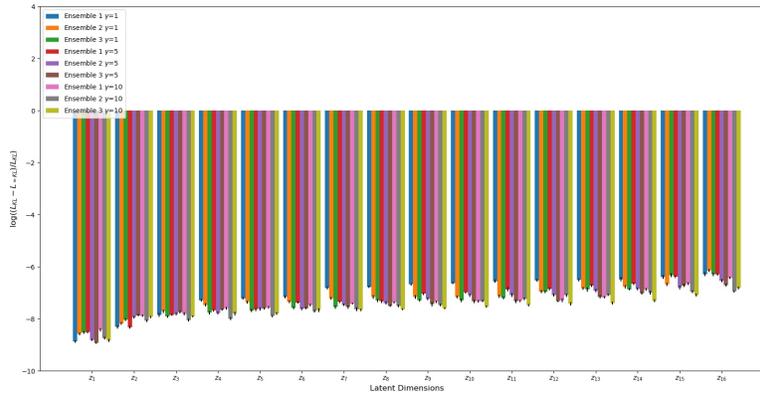


(c) VAE_E2 with latent dimension of 10, CelebA data.

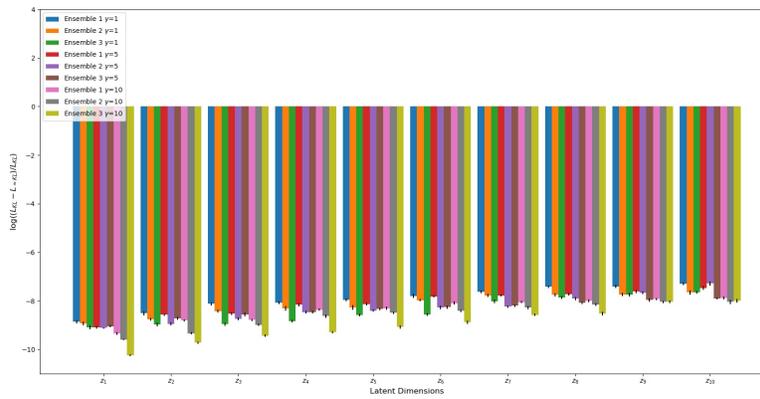
Figure 9: The “polarized regime” comparison between models in VAE_E2. The results are over 5 runs of training on the *CelebA* dataset.



(a) VAE_E₃ with latent dimension of 32, CelebA data.

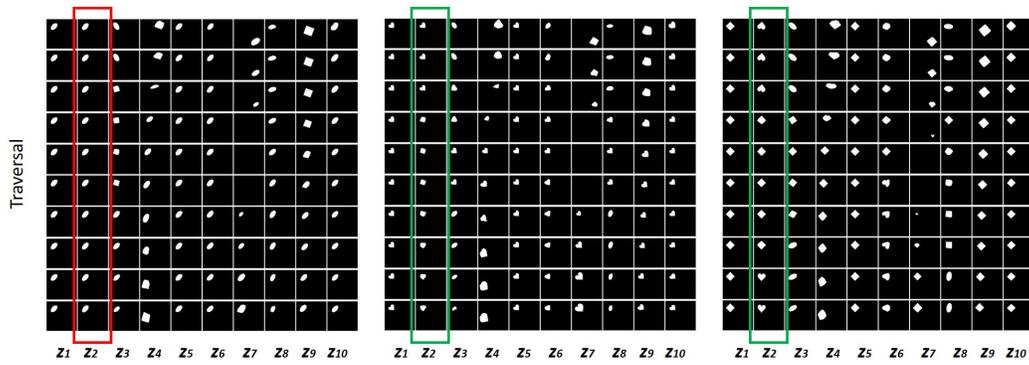


(b) VAE_E₃ with latent dimension of 16, CelebA data.

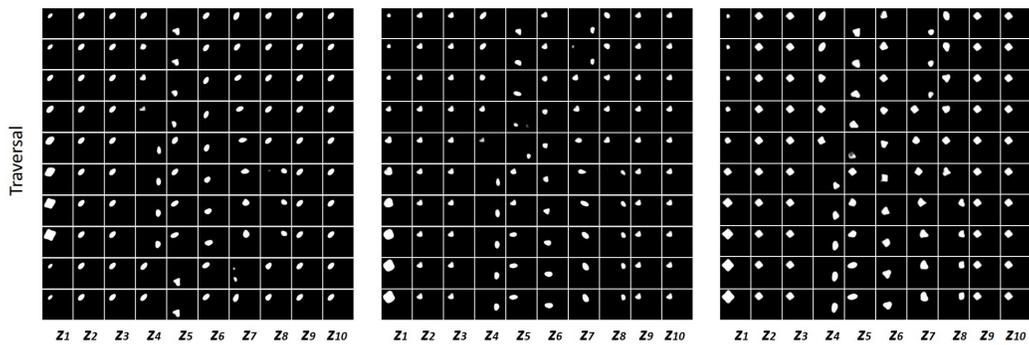


(c) VAE_E₃ with latent dimension of 10, CelebA data.

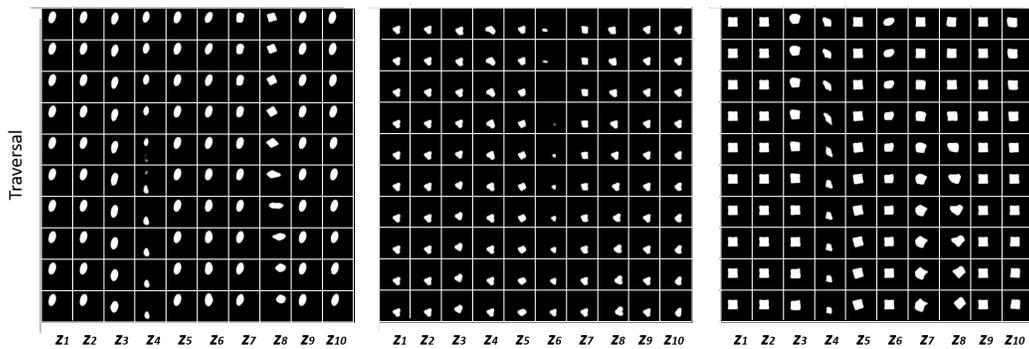
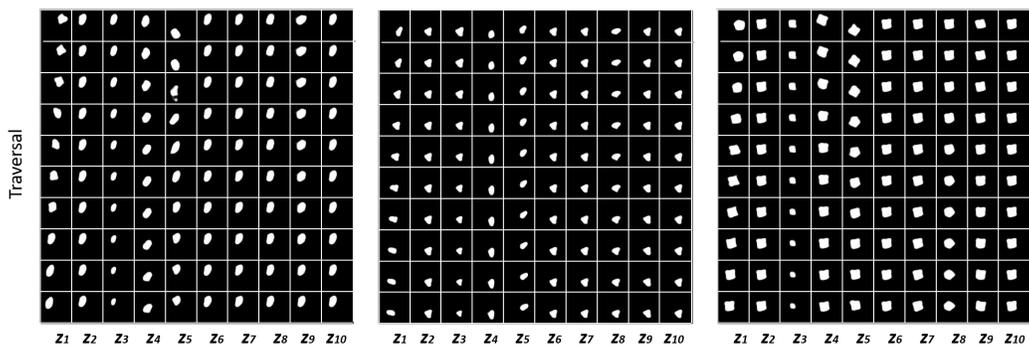
Figure 10: The “polarized regime” comparison between models in the VAE_E₃ on the *CelebA* dataset.



(a) VAE_E2

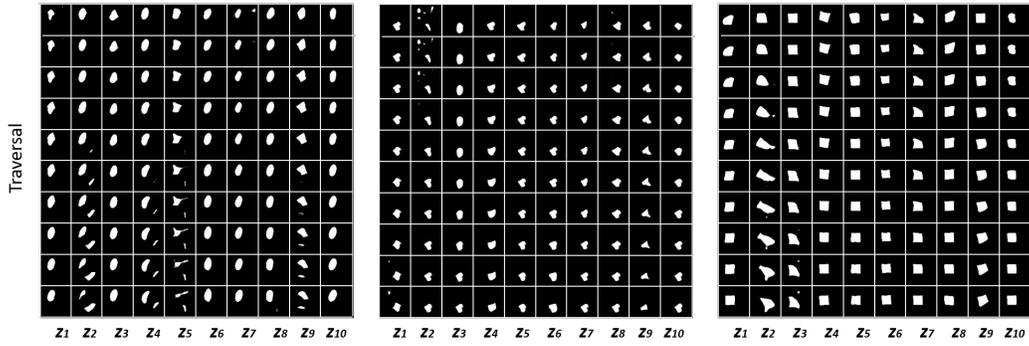


(b) Single VAE

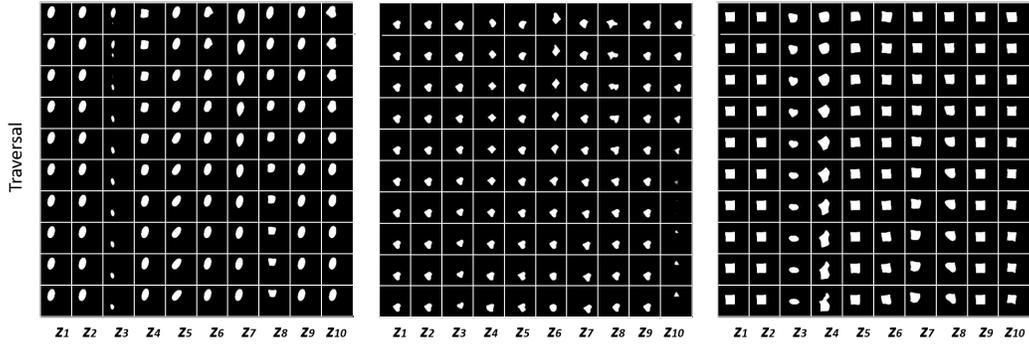
(c) β -VAE

(d) FactorVAE

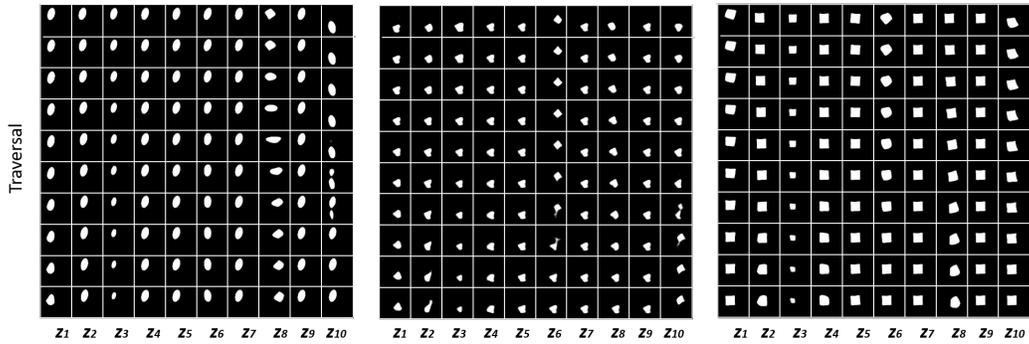
Figure 11: Latent traversal on three different input images using VAE_E2, a single VAE and the state-of-the-art VAE models with 10 dimensional latent representation. The three input images are ellipse, heart and square shapes as shown in the last column.



(e) DIP-VAE-I



(f) DIP-VAE-II



(g) TC-VAE

Figure 11: (cont.) Latent traversal on three different input images using VAE_{E₂}, a single VAE and the state-of-the-art VAE models with 10 dimensional latent representation. The three input images are ellipse, heart and square shapes as shown in the last column.



Figure 12: Latent traversal on two different input images of *CelebA* dataset using VAE $_E_2$ with latent dimension of 16.

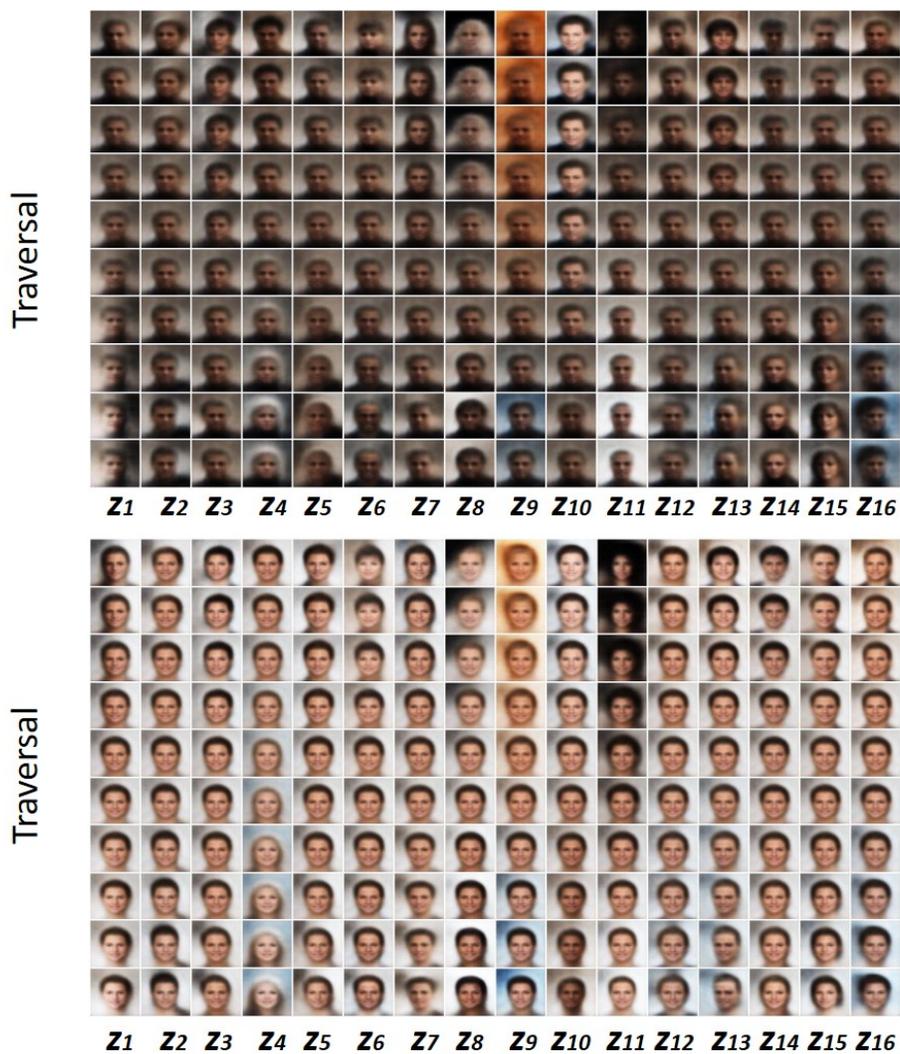


Figure 13: Latent traversal on two different input images of *CelebA* dataset using a single VAE with latent dimension of 16.