

---

# GroupifyVAE: from Group-based Definition to VAE-based Unsupervised Representation Disentanglement

## Appendix

---

### A. Preliminary

#### A.1. Group Theory Basics

**Group:** A set  $G$  together with a binary operation  $\circ$  as:  $G \times G \rightarrow G$  satisfying the following conditions:

- **Associativity:**  $\forall a, b, c \in G, s.t. (a \circ b) \circ c = a \circ (b \circ c)$
- **Identity:**  $\exists e \in G, s.t. \forall a \in G, e \circ a = a \circ e = a$
- **Inverse:**  $\forall a \in G, \exists a^{-1} \in G : a \circ a^{-1} = a^{-1} \circ a = e$

It is customary to represent a Group with a set  $G$  and the binary operation  $\circ$  pair  $(G, \circ)$ . When the binary operation is clear, it is customary to represent Group  $(G, \circ)$  as  $G$ . In this case, it is customary to use multiplication to represent the binary operation  $\circ$ , i.e.,  $a \circ b = ab, \forall a, b \in G$ .

**Abelian Group:** A group  $(G, \circ)$  satisfies the commutative condition:  $\forall a, b \in G, a \circ b = b \circ a$

**Subgroup:** A Group  $(G, \circ)$ ,  $A$  is a subset of  $G$ , if  $(A, \circ)$  forms a Group, then  $A$  is the subgroup of  $G$ .

**Equivalence Relation:**  $A$  is a set, the subset  $R$  of set  $A \times A = \{(a, b); \forall a, b \in A\}$  is a **Relation** on set  $A$ , and set  $R$  satisfying the following conditions: (It is customary represent  $(a, b) \in R$  as  $a \sim b$ )

- **Reflexive:**  $\forall a \in A, s.t. a \sim a$
- **Symmetric:** if  $a \sim b$ , then  $b \sim a$
- **Transitive:** if  $a \sim b$  and  $b \sim c$ , then  $a \sim c$

$\forall a \in A$ , we represent all the elements in  $A$  that equivalent to  $a$  as  $\bar{a}$ , i.e.,  $\bar{a} = \{b \in A; b \sim a\}$ ,  $\bar{a}$  is **Equivalent Class**.

**Isomorphism:**  $(G, \cdot), (G', \circ)$  are two Groups, mapping  $f : G \rightarrow G'$  is a **Homomorphism** between  $G$  and  $G'$ , if  $\forall a, b \in G, f(a \cdot b) = f(a) \circ f(b)$ , if mapping  $f$  is a bijection, then mapping  $f$  is a Isomorphism between  $G$  and  $G'$ .

#### Related Examples of Groups:

(i) **Multiplicative group of integers modulo  $n$ :**  $n$  is a positive integer, we define a Equivalence Relation on  $\mathbb{Z}$  as  $a \sim b \Leftrightarrow n|a - b$ , i.e.,  $a = b(mod n)$ . The relation divide  $\mathbb{Z}$  into  $n$  Equivalent Classes:  $\overline{0}, \overline{1}, \dots, \overline{n-1}$ , where  $\bar{i}$  represents the Equivalent Class containing  $i$ , i.e.,  $\bar{i} = \{m \in \mathbb{Z} | m = i(mod n)\}$ , let  $Z_n = \{\overline{0}, \overline{1}, \dots, \overline{n-1}\}$ ,  $Z_n$  with binary operation:  $\bar{a} + \bar{b} = \overline{a+b}$  forms a Group, it is denoted by  $\mathbb{Z}/n\mathbb{Z}$ .

(ii)  **$n$ -th root unity Group:**  $n$  is a positive integer.  $C_n = \{e^{\frac{2\pi i a}{n}} | 0 \leq a \leq n-1\}$ , then  $C_n$  with complex multiplication forms a Group. For mapping:  $f : C_n \rightarrow (Z_n, +), e^{\frac{2\pi i a}{n}} \mapsto \bar{a}$ , then

$$f\left(e^{\frac{2\pi i a}{n}} \cdot e^{\frac{2\pi i b}{n}}\right) = f\left(e^{\frac{2\pi i a+b}{n}}\right) = \overline{a+b} = \bar{a} + \bar{b} = f\left(e^{\frac{2\pi i a}{n}}\right) \cdot f\left(e^{\frac{2\pi i b}{n}}\right) \quad (1)$$

Therefore,  $f$  is a homomorphism,  $f$  is a bijection, and  $f$  is a isomorphism. In this paper, we designed two constrains to let the encoder and decoder forms a isomorphism mapping  $f$ .

**Congruence Class:**  $n$  is a positive integer, we define a Equivalence Relation on  $\mathbb{Z}$  as  $a \sim b \Leftrightarrow n|a - b$ , i.e.,  $a = b(mod n)$ . The relation divide  $\mathbb{Z}$  into  $n$  Equivalent Classes:  $\overline{0}, \overline{1}, \dots, \overline{n-1}$ , which are congruence classes and also the elements of the multiplicative group of integers modulo  $n$ .

**Subgroup Generated by A:** If  $A$  is a nonempty subset of the Group  $G$ , the **Set Generated by A**, denoted by  $\langle A \rangle$ , is the set defined by  $\langle A \rangle = \{a \in G | a = a_1 a_2 \dots a_n \text{ with either } a_i \in A \text{ or } a_i^{-1} \in A\}$ , then the set  $\langle A \rangle$  is a subgroup of  $G$ .

**Symmetry Group and Permutation Group:** For a nonempty set  $\Sigma$ , the one-to-one mappings on itself  $\sigma$  called a **Permutation**,  $S(\Sigma)$  denotes the set containing all the Permutations on  $\Sigma$ ,  $S(\Sigma)$  with mapping compound as binary operator forms a Group called Symmetry Group, the subgroup of  $S(\Sigma)$  called Permutation Group. In particular, Symmetry Group of set  $\{1, 2, \dots, n\}$  is denoted by  $S_n$ .

**Group Action:**  $\Sigma$  is a set,  $S(\Sigma)$  is the Symmetry Group on  $\Sigma$ , homomorphism  $\forall f : G \rightarrow S(\Sigma)$  called **permutation representation** on set  $\Sigma$ , then  $f(g)$  is a Permutation on  $\Sigma$  called  $g$  action on  $\Sigma$ .  $\forall a \in \Sigma$ , define  $ga = f(g)a$

**Symmetry:** A symmetry of a geometric figure is a rearrangement of the figure preserving the arrangement of its sides and vertices as well as its distances and angles.

Exmaple of symmetry: Take rectangle as an example in Figure 1, it is easy to see that a rotation of  $180^\circ$  or  $360^\circ$  are symmetries for returning a rectangle with the same orientation as the original one and the same relationship among the vertices. Similarly, A reflection either the vertical axis or the horizontal axis can also be seen to be a symmetry.

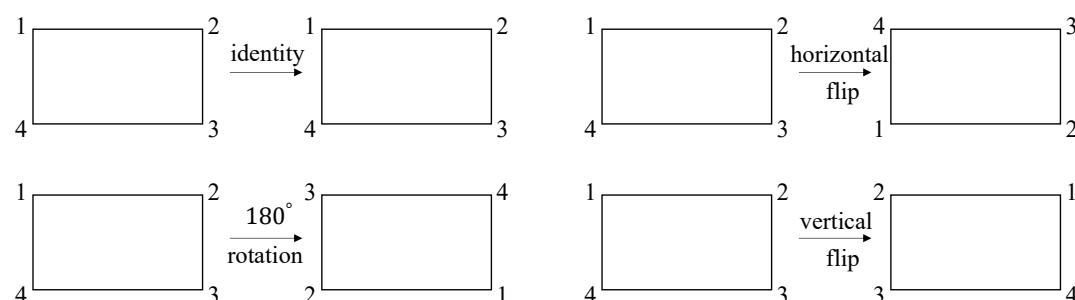


Figure 1. Rigid motions of a rectangle

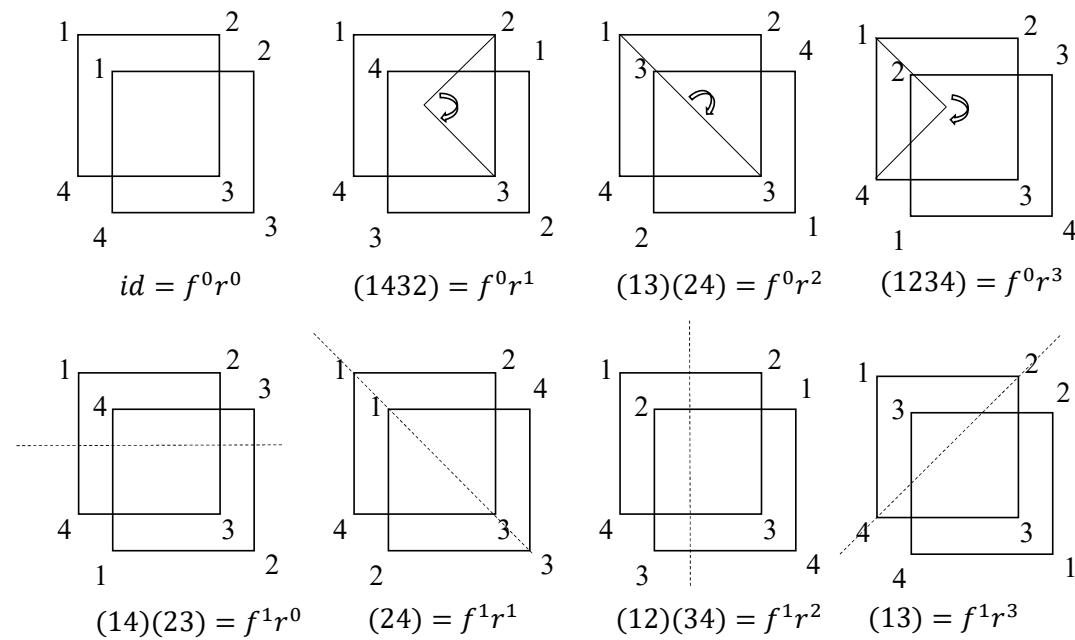


Figure 2. 4-th Dihedral Group ( $D_4$ ): The groups of symmetries for square.

**Rigid motion:** A map from the plane to itself preserving the symmetry of an object is called a rigid motion. In particular, rigid motions of a regular  $n$ -sided polygon are permutations of the vertices.

**$n$ th Dihedral Group:** The group of rigid motions of a regular  $n$ -sided polygon, which is denoted by  $D_n$ ,  $D_n \subseteq S_n$ .

110  **$D_4$  as an example:** The vertices of a square are numbered by  $\{1, 2, 3, 4\}$ , which is in analog with image dataset, the rigid  
 111 motions is in analog with actions in latent space. We often abbreviate permutation  $1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4, 4 \rightarrow 1$  as  $(1234)$ .  
 112 The elements of Group  $D_4$  are shown in Figure 2, from which we know that all of the transformations are compounded by  
 113 basic permutation horizontal flip  $f$  and rotate 90 degrees clockwise  $r$ . Please note that  $f$  and  $r$  are in analog with disentangled  
 114 factors. What's more, The Group can be generated by these basic permutations:  $D_4 = \langle f, r | f^2 = 1, r^4 = 1, fr = r^{-1}f \rangle$ .  
 115 The constraints  $f^2 = 1, r^4 = 1, fr = r^{-1}f$  is in analog with the group constraints in this paper.  
 116

## 117 B. Proof for Theorem 1

119 *Proof.*  $\Rightarrow$ ) First, we show that when the existance of isomorphism  $G \sim \Phi = \langle \varphi_1, \varphi_2, \dots, \varphi_m \rangle$  and  $\varphi_i(f(w)) =$   
 120  $f(\overline{w + c_i})$  will make the mapping  $c$  be equivariant between the actions on  $W$  and  $Z$ :  $g \cdot c(w) = c(g \cdot w)$ . Then, we proof  
 121 that each  $Z_i$  is affected only by  $G_i$ .

122 Take  $\forall g \in G, G = \mathbb{Z}/n\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z} \times \dots \times \mathbb{Z}/n\mathbb{Z}$ , here we assume  $g = \sum_i k_i e_i$ . In this setting, the group action on  $W$  is:  
 123  $g \cdot w = w + \sum_i k_i e_i, \forall w \in W$ . Similarly, the group actions on  $Z$  is:  $g \cdot z = z + \sum_i k_i e_i, \forall z \in Z$ .  
 124

125 For the group  $\Phi = \langle \varphi_1, \varphi_2, \dots, \varphi_m \rangle$ , where  $\varphi_i \in S\{f(w)\}$ ,  $S\{f(w)\}$  is symmetry group of  $\{f(w)\}$ , the definition of  
 126  $\varphi_i$  is:

$$127 \quad \varphi_i(f(w)) = b^{-1}(\overline{b(f(w)) + e_i}) = b^{-1}(\overline{c(w) + e_i}), \quad (2)$$

128 where  $b^{-1} = d$  denotes the decoder, which can be regard as the inversion of encoder  $b$ . And the mapping  $c$  is defined as  
 129  $c = b \circ f$ . Assume there is a isomorphism  $\tau : G \rightarrow \Phi$ , we have  $\text{ker}(\tau) = \{e_G\}$ ,  $e_G$  is the identity element in group  $G$ , and  
 130  $e_\Phi$  is identity element in group  $\Phi$ , so that  $\tau(e_G) = e_\Phi$ . Let  $U = \{e_i \in G | 0 \leq i \leq m\}$ , where  $e_i$  is the identity element of  
 131 dimension  $i$  in  $G$ .  
 132

133 Assume  $\exists i, 0 \leq i \leq m$ , s.t.  $\tau^{-1}(\varphi_i) = g, \forall g \in G$  but  $g \notin U$ , since  $\tau$  is a bijection,

$$134 \quad |\Phi| = |\langle \varphi_1, \varphi_2, \dots, \varphi_m \rangle| = |\langle e_1, e_2, \dots, e_{i-1}, g, e_{i+1}, \dots, e_m \rangle| < |G|. \quad (3)$$

137 However, since  $\tau$  is a isomorphism, we have  $|\Phi| = |G|$ , which is contradicts the above equation. Therefore,  $\forall 0 \leq i \leq m, \tau^{-1}(\varphi_i) \in U, \tau(e_i) = \varphi_j$ , by rearrange subscript of  $\varphi_i$ . May as well set  $\tau(e_i) = \varphi_i$

139 Therefore, n times compose  $\varphi_i$ , since  $\tau$  is a isomorphism, we have

$$141 \quad \varphi_i^n = \tau(e_i)^n = \tau(ne_i) = \tau(e_G) = e_\Phi = id. \quad (4)$$

143 Recall that another condition:  $\varphi_i(f(w)) = f(\overline{w + c_i})$ , this condition with n times compose:

$$144 \quad \varphi_i^n(f(w)) = f(\overline{w + nc_i}) = f(w), \quad (5)$$

146 which indicate that  $c_i = (k+1)e_i$ , if  $k \neq 0$ , then

$$148 \quad \text{Order}(\varphi_i) = \frac{[n, k+1]}{k+1} < n, \quad (6)$$

150 which contradict with  $\text{order}(\varphi_i) = n$ . Therefore  $k = 0$  (Please note that we do not consider the case of mutually exclusive  
 151 Euler groups, if they are mutually exclusive, after redefining the generator, change the scale, this group is remain the same).  
 152 Therefore in this case,

$$154 \quad \varphi_i(f(w)) = f(\overline{w + e_i}), \Phi = \{\varphi_1^{k_1} \varphi_2^{k_2} \dots \varphi_n^{k_n} | k_i \in Z, \overline{k_i} = k_i \mod n\}. \quad (7)$$

156 We have

$$157 \quad \forall g \in G, b^{-1}(g \cdot c(w)) = b^{-1}(\overline{\sum_i k_i e_i + c(w)}) = \varphi_1^{k_1} \varphi_2^{k_2} \dots \varphi_n^{k_n}(f(w)) = f(\overline{w + \sum_i k_i e_i}) = f(g \cdot w), \quad (8)$$

160 then we have  $g \cdot c(w) = b \circ f(g \cdot w) = c(g \cdot w)$ , i.e. The map  $c$  is equivariant between the actions on  $W$  and  $Z$ .

161 Then we proof that each  $Z_i$  is affected only by  $G_i$ .  $\forall g_i \in G_i$ , for group action on  $f(w)$

$$163 \quad g_i \cdot f(w) = \varphi_i^{k_i}(f(w)) = f(\overline{w + ke_i}) = b^{-1}(\overline{ke_i + c(w)}). \quad (9)$$

165 Based on the above equation, we have

$$g_i \cdot c(w) = b \circ b^{-1}(\overline{ke_i + c(w)}) = \overline{ke_i + c(w)}. \quad (10)$$

166 Therefore,  $Z_i$  is affected only by  $g_i$  for the decomposition  $Z = Z_1 \times Z_2 \times \dots \times Z_n$ . Summarizing, the representation  $Z$  is  
167 disentangled with respect to  $G$ .

170  $\Leftrightarrow$ ) Review that definition of  $\varphi_i$  is

$$\varphi_i(f(w)) = b^{-1}(\overline{b(f(w)) + e_i}) = b^{-1}(\overline{c(w) + e_i}), \quad (11)$$

174 the mapping  $c$  is equivariant between the actions on  $W$  and  $Z$ :  $g \cdot c(w) = c(g \cdot w)$ , we have

$$\varphi_i(f(w)) = b^{-1}(\overline{c(w + e_i)}). \quad (12)$$

177 Therefore,  $\varphi_i(f(w)) = f(\overline{w + e_i})$ ,  $c = b \circ f$  is satisfied, and we define the mapping  $\tau$  as:

$$\tau(\varphi) = g : \varphi(f(w)) = b^{-1}(\overline{c(w + g)}). \quad (13)$$

181 Because

$$\varphi_s(\varphi_t(f(w))) = \varphi_s(b^{-1}(\overline{c(w + g_s)})) = \varphi_s(f(w + g_s)) = b^{-1}(\overline{c(w) + g_s + g_t}), \quad (14)$$

183 and we have

$$\tau(\varphi_s \varphi_t) = g_s + g_t = \tau(\varphi_s) + \tau(\varphi_t). \quad (15)$$

185 Consequently,  $\tau$  is a homomorphism. Due to  $e_\Phi(f(w)) = f(w) = b^{-1}(c(w))$ ,  $\tau(e_\Phi) = e_G$ . Assume that there exist  
186  $\varphi_e \neq e_\Phi$ , s.t.  $\tau(\varphi_e) = e_G$ , we have

$$\varphi_e(f(w)) = b^{-1}(c(w)) = f(w). \quad (16)$$

189 Therefore,  $\varphi_e = id = e_\Phi$ , and  $\tau^{-1}(e_G) = \{e_\Phi\}$ ,  $\tau$  is injective. We take  $\forall g \in G, g = (\overline{k_1}, \overline{k_2}, \dots, \overline{k_m})$ , then

$$b^{-1}(c(w + g)) = b^{-1}(c(w + \overline{k_1}e_1 + \overline{k_2}e_2 + \dots + \overline{k_m}e_m)) = \dots = \varphi_1 \circ \dots \circ \varphi_m(f(w)). \quad (17)$$

192 Moreover,  $\varphi_1 \circ \dots \circ \varphi_m \in \Phi$ ,  $\tau$  is surjective, and  $\tau$  is a isomorphism.

193 Q.E.D.

## 195 C. Proof for Theorem 2

196 *Proof.* The theorem is proofed in the following two steps, In the first step, we prove that the set generated by the  
197 transformation  $\Phi = \langle \varphi_i | \varphi_i \varphi_j = \varphi_j \varphi_i, \varphi_i^n = e_\Phi = id, \forall 0 \leq i, j \leq m \rangle$  forms a group under the compound operation.  
198 In the second step, we prove that the necessary and sufficient condition for the existence of isomorphism is that the Abel  
199 constraint and the Order constraint are satisfied at the same time.

200 **Step 1:** To verify that a set forms a group under a certain operation, we only need to verify that the elements in the set  
201 satisfied the following three requirements: 1. Associativity 2. Identity 3. Inverse. We take  $\forall \varphi_i \in S\{f(w)\}$ ,  $S\{f(w)\}$  is the  
202 symmetry group on images set  $\{f(w)\}$ . Therefore  $\varphi_i$  is a mapping:  $\varphi_i : X \rightarrow X$ , then  $\forall \varphi_i, \varphi_j, \varphi_k \in \{\varphi_i, 0 \leq i \leq m\}$ ,

$$(\varphi_i \varphi_j) \varphi_k = \varphi_i(\varphi_j \varphi_k), \quad (18)$$

203 thus generators of  $\Phi$  satisfy Associativity. However, the dose the element of  $\Phi$  satisfy Associativity is unknown, we take  
204  $\forall \varphi_s, \varphi_t, \varphi_l \in \Phi$ ,

$$(\varphi_t \varphi_s) \varphi_l = \left( \prod_{s_i} \overline{\varphi_{s_i}^{k_{s_i}}} \prod_{t_i} \overline{\varphi_{t_i}^{k_{t_i}}} \right) \prod_{l_i} \overline{\varphi_{l_i}^{k_{l_i}}} = \prod_{s_i} \overline{\varphi_{s_i}^{k_{s_i}}} \left( \prod_{t_i} \overline{\varphi_{t_i}^{k_{t_i}}} \prod_{l_i} \overline{\varphi_{l_i}^{k_{l_i}}} \right) = \varphi_t(\varphi_s \varphi_l). \quad (19)$$

205 Therefore, the mapping set  $\Phi$  satisfy Associativity under the compound operation. In the following part, we will find the unit  
206 element and inverse element of  $\Phi$ , for the generator elements,  $\forall 0 \leq i, j \leq m$ , we have  $\varphi_i^n \varphi_j = id \cdot \varphi_j = \varphi_j$ , Therefore,  
207 for general elements we take  $\forall \varphi \in \Phi, \varphi = \prod_i \varphi_i^{k_i}$ , we have:

$$\varphi_i^n \varphi = (id \cdot \varphi_j) \varphi_j^{\overline{k_j - 1}} \prod_{i \neq j} \overline{\varphi_i^{k_i}} = \varphi_j^{\overline{k_j}} \prod_{i \neq j} \overline{\varphi_i^{k_i}} = \prod_i \overline{\varphi_i^{k_i}} = \varphi, \quad (20)$$

and this implies that we have unit  $e = \varphi_i^n = id$  in  $\Phi$ . Identity is satisfied. For the generator element,  $\forall 0 \leq i \leq m$ , we have  $\varphi_i^{\overline{k}} \varphi_i^{\overline{n-k}} = \varphi_i^{\overline{n}} = id = e$ ,  $0 \leq k \leq n$ , for general element,  $\forall \varphi \in \Phi, \varphi = \prod_i \varphi_i^{\overline{k_i}}$ , we have:

$$\varphi \prod_i \varphi_i^{\overline{n-k_i}} = \prod_{i \{i \neq j\}} \varphi_i^{\overline{k_i}} (\varphi_j^{\overline{k_j}} \varphi_j^{\overline{n-k_j}}) \prod_{i \{i \neq j\}} \varphi_i^{\overline{n-k_i}} = \dots = id, \quad (21)$$

which implies that for general element  $\forall \varphi \in \Phi, \varphi = \prod_i \varphi_i^{\overline{k_i}}$ , we have the inverse element  $\varphi^{-1} = \prod_i \varphi_i^{\overline{n-k_i}}$  in  $\Phi$ . Inverse is satisfied.

Summarizing,  $\Phi$  is a group. Moreover, for generative element  $\varphi_i, \varphi_j, 0 \leq i, j \leq m$ , we have  $\varphi_i \varphi_j = \varphi_j \varphi_i$ , for general element  $\varphi_s, \varphi_t \in \Phi$ ,

$$\varphi_t \varphi_s = \prod_{s_i} \varphi_{s_i}^{\overline{k_{s_i}}} \prod_{t_i} \varphi_{t_i}^{\overline{k_{t_i}}} = \prod_{t_i} \varphi_{t_i}^{\overline{k_{t_i}}} \prod_{s_i} \varphi_{s_i}^{\overline{k_{s_i}}} = \varphi_s \varphi_t. \quad (22)$$

By pairwise exchange of generators, the  $\varphi_{s_i}^{\overline{k_{s_i}}}$  is moved to the head of the equation one by one, then, equation 22 is obtained.

Summarizing,  $\Phi$  is a Abelian group.

Q.E.D.

**Step 2:** In this step, we prove that the necessary and sufficient condition for the existence of isomorphism  $\mathbb{Z}/n\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z} \times \dots \mathbb{Z}/n\mathbb{Z} \sim < \varphi_1, \varphi_2, \dots, \varphi_m >$  is that  $0 \leq i, j \leq m, \varphi_i \varphi_j = \varphi_j \varphi_i$ , and  $0 \leq i \leq m, \varphi_i^n = e$  are satisfied simultaneously.

$\Rightarrow$ ) Assume the isomorphism is  $h : G = \mathbb{Z}/n\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z} \times \dots \mathbb{Z}/n\mathbb{Z} \rightarrow \Phi = < \varphi_1, \varphi_2, \dots, \varphi_m >$ , for the generator elements,  $\forall \varphi_i, \varphi_j \in \{\varphi_i | 0 \leq i \leq m\}$ , there are  $g_s = \sum_{s_i} k_{s_i} e_{s_i}, g_t = \sum_{t_i} k_{t_i} e_{t_i} \in G$ , s.t.  $h(g_s) = \varphi_i, h(g_t) = \varphi_j$ , since  $G$  is an Abelian group.

$$\varphi_i \varphi_j = h(g_s)h(g_t) = h(g_s g_t) = h(g_t g_s) = h(g_t)h(g_s) = \varphi_j \varphi_i. \quad (23)$$

For the generator elements in  $\Phi$ ,  $0 \leq i \leq m, \varphi_i$ , we have  $g_l = \sum_{l_i} k_{l_i} e_{l_i} \in G$ , s.t.  $h(g_l) = \varphi_i$ , compose itself n times,

$$\varphi_i^n = h(g_l)^n = h(n g_l) = h(\sum_{l_i} k_{l_i} n e_{l_i}) = h(\sum_{l_i} k_{l_i} e_G) = h(e_G) = e_\Phi. \quad (24)$$

In equation 24,  $h(e_G) = e_\Phi$  is hold for  $h$  is a isomorphism, and  $ker(h) = e_G$ . Sufficiency is proven.

$\Leftarrow$ ) In the following part, we prove that when the two conditions are satisfied at the same time, the mapping  $\tau$  we found is an isomorphism. Considering mapping  $\tau : \Phi \rightarrow G$ , The one-hot vector  $(0, \dots, \overline{1}, \dots, 0) \in G$ , that position i is 1 else is 0, is denoted by  $e_i$ . The definition of  $\tau$ ,

$$\begin{cases} \tau : \varphi_i \mapsto e_i; \\ \tau : \varphi_i \varphi_j \mapsto e_i + e_j. \end{cases} \quad (25)$$

For general elements in  $\Phi$ ,  $\forall \phi_t, \phi_s \in \Phi$ , where  $\phi_t = \prod_{i_t} \varphi_{i_t}^{\overline{k_{i_t}}}, \phi_s = \prod_{i_s} \varphi_{i_s}^{\overline{k_{i_s}}}$ , Note that  $\phi_s$  is written in this form only if the set  $\Phi$  is proved to be a group, which is hold by the two conditions.

$$\tau(\phi_t \phi_s) = \tau \left( \prod_{i_t} \varphi_{i_t}^{\overline{k_{i_t}}} \prod_{i_s} \varphi_{i_s}^{\overline{k_{i_s}}} \right) = \tau \left( \prod_{i_l} \varphi_{i_l}^{\overline{k_{i_l}}} \right) = \sum_{i_l} \overline{k_{i_l}} e_{i_l}, \quad (26)$$

and please note that the subscript can be merged since  $\varphi_i \varphi_j = \varphi_j \varphi_i$ , where  $i_l$  denotes the subscript after  $i_t, i_s$  merged. The summation on the right side of equation 26 is devided into two parts,

$$\sum_{i_l} \overline{k_{i_l}} e_{i_l} = \sum_{i_t} \overline{k_{i_t}} e_{i_t} + \sum_{i_s} \overline{k_{i_s}} e_{i_s} = \tau \left( \prod_{i_l} \varphi_{i_l}^{\overline{k_{i_l}}} \right) + \tau \left( \prod_{i_l} \varphi_{i_l}^{\overline{k_{i_l}}} \right) = \tau(\phi_t) + \tau(\phi_s). \quad (27)$$

Consequently,  $\tau(\phi_t \phi_s) = \tau(\phi_t) + \tau(\phi_s)$ , this implies that  $\tau$  is a homomorphism. Then, only need to prove that the mapping  $\tau$  is bijective. First we prove  $\tau$  is injective, i.e.  $ker(\tau) = \{e_G\}$ , after that we prove  $\tau$  is surjective, i.e. Find the inverse image of any element.

275 It's not hard to obtain  $\tau(\phi) = \tau(e_\Phi \cdot \phi) = \tau(e_\Phi) + \tau(\phi)$ , therefore  $\tau(e_\Phi) = (0, \dots, 0) = e_G$ , i.e.  $e_\Phi \in \tau^{-1}(e_G)$ , assume  
 276 there is  $\phi_e = \prod_{i_e} \varphi_{i_e}^{\overline{k_{i_e}}}$ , s.t.  $\tau(\phi_e) = e_G$   
 277

$$\begin{aligned} \phi_e &= \prod_{i_e} \varphi_{i_e}^{\overline{k_{i_e}}}, \tau(\phi_e) = \sum_{i_e} k_{i_e} e_{i_e} = e_G \Rightarrow k_{i_e} = l \cdot n; \\ \phi_e &= \prod_{i_e} \varphi_{i_e}^{\overline{k_{i_e}}} = \prod_{i_e} e_\Phi = e_\Phi, \end{aligned} \quad (28)$$

283 which means,  $\tau^{-1}(e_G) = \{e_\Phi\}$ , mapping  $\tau$  is injective. In the following part, only need to prove that for any element in  $G$ ,  
 284 we can find the inverse image of it. We take  $\forall g \in G, g = (\overline{k_1}, \overline{k_1}, \dots, \overline{k_m}) \in \mathbb{Z}/n\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z} \times \dots \times \mathbb{Z}/n\mathbb{Z}$ ,  
 285

$$g = (\overline{k_1}, \dots, 0) + \dots + (0, \dots, \overline{k_m}) = \sum_i \overline{k_i} e_i = \tau \left( \prod_i \varphi_i^{\overline{k_i}} \right). \quad (29)$$

290 Hence,  $\tau^{-1}(g) = \prod_i \varphi_i^{\overline{k_i}} \in \Phi$ , This implies that mapping  $\tau$  is surjective. Since mapping  $\tau$  is bijective and homomorphism,  
 291 we conclude that  $\tau$  is a isomorphism.  
 292

293 Q.E.D.

## D. Proof for Theorem 3

294 *Proof.* Here we first introduce image space  $X = \{x_1, x_2, \dots, x_N\}$  and representation space  $Z = \{z_1, z_2, \dots, z_N\}$ , the  
 295 encoder  $b : X \rightarrow Z$ , and the decoder  $d : Z \rightarrow X$ , the definition  $\varphi_i : X \rightarrow X, 0 \leq i \leq m$  is  $\varphi_i(x) = d(a_i(\overline{b(x)}))$ , where  $a_i$   
 296 is the group action of  $e_i$  on  $Z$ .

300  $\Rightarrow$ ) Obviously, we can obtain the expend form of  $\varphi_j \circ \varphi_i$  from the above definition,

$$(\varphi_j \circ \varphi_i)(x) = \varphi_j(\varphi_i(x)) = \varphi_j(d(a_i(\overline{b(x)}))) = d(a_j(\overline{b(d(a_i(\overline{b(x)}))))}). \quad (30)$$

304 For the condition:  $\forall \varphi_i, \varphi_j \in \Phi, 0 \leq i, j \leq m$ , we have  $\varphi_i \varphi_j = \varphi_j \varphi_i$ . The constrain on specific image can be obtained by  
 305

$$\varphi_i(\varphi_j(x)) = \varphi_j(\varphi_i(x)) \Rightarrow \varphi_i(\varphi_j(x)) - \varphi_j(\varphi_i(x)) = 0. \quad (31)$$

307 For the set of combinations of factors  $C = \{(i, j) | \forall i, j \in I\}$ , and the set containing images  $X$ , we have  
 308

$$\mathcal{L}_a = \sum_{x \in X} \sum_{(i,j) \in C} \|\varphi_i(\varphi_j(x)) - \varphi_j(\varphi_i(x))\| \text{ is optimized.} \quad (32)$$

312 The Abel loss is obtained. For the Order loss, we first consider  $n$  times compose of the same mapping  $\varphi_i$ ,

$$\varphi_i^n = e_\Phi = id, \Rightarrow \varphi_i^{n-1} \varphi_i = \varphi_i \varphi_i^{n-1} = e_\Phi = id. \quad (33)$$

316 Therefore,  $\varphi_i^{-1} = \varphi_i^{n-1}$ . for a single image  $x$ , we have

$$\varphi_i(\varphi_i^{n-1}(x)) = \varphi_i(\varphi_i^{-1}(x)) = x \Rightarrow \varphi_i(\varphi_i^{-1}(x)) - x = 0. \quad (34)$$

319 Thus, for the set of factors  $I$  and the set containing all images  $X$ , we have  
 320

$$\sum_{x \in X} \sum_{i \in I} \|\varphi_i(\varphi_i^{-1}(x)) - x\| \text{ is optimized.} \quad (35)$$

324 For eliminating the bias of optimization, we optimize a symmetry form of the Order Loss, and we have  
 325

$$\mathcal{L}_o = \sum_{x \in X} \sum_{i \in I} (\|\varphi_i(\varphi_i^{-1}(x)) - x\| + \|\varphi_i^{-1}(\varphi_i(x)) - x\|) \text{ is optimized.} \quad (36)$$

328 The Order Loss is obtained.  
 329

330  $\Leftrightarrow$ ) When the Abel Loss  $\mathcal{L}_a$  is optimized,  $\forall x \in X$ , we have

$$331 \quad 332 \quad 333 \quad 334 \quad 335 \quad 336 \quad 337 \quad 338 \quad 339 \quad 340 \quad 341 \quad 342 \quad 343 \quad 344 \quad 345 \quad 346 \quad 347 \quad 348 \quad 349 \quad 350 \quad 351 \quad 352 \quad 353 \quad 354 \quad 355 \quad 356 \quad 357 \quad 358 \quad 359 \quad 360 \quad 361 \quad 362 \quad 363 \quad 364 \quad 365 \quad 366 \quad 367 \quad 368 \quad 369 \quad 370 \quad 371 \quad 372 \quad 373 \quad 374 \quad 375 \quad 376 \quad 377 \quad 378 \quad 379 \quad 380 \quad 381 \quad 382 \quad 383 \quad 384$$

$$\varphi_i(\varphi_j(x)) - \varphi_j(\varphi_i(x)) = 0 \Rightarrow \varphi_i(\varphi(x)) = \varphi_j(\varphi_i(x)). \quad (37)$$

Therefore,  $\forall \varphi_i, \varphi_j \in \Phi, 0 \leq i, j \leq m$ , we have  $\varphi_i \varphi_j = \varphi_j \varphi_i$ , we obtain the Abel constraint. When the Order Loss  $\mathcal{L}_o$  is optimized,  $\forall x \in X$ , we have

$$\varphi_i(\varphi_i^{-1}(x)) - x = 0 \Rightarrow \varphi_i(\varphi_i^{-1}(x)) = x. \quad (38)$$

This implies that

$$\varphi_i^n = \varphi_i \circ \varphi_i^{n-1} = \varphi_i \circ \varphi_i^{-1} = e. \quad (39)$$

Therefore, we obtain the Order constraint. The Group constraints are satisfied.

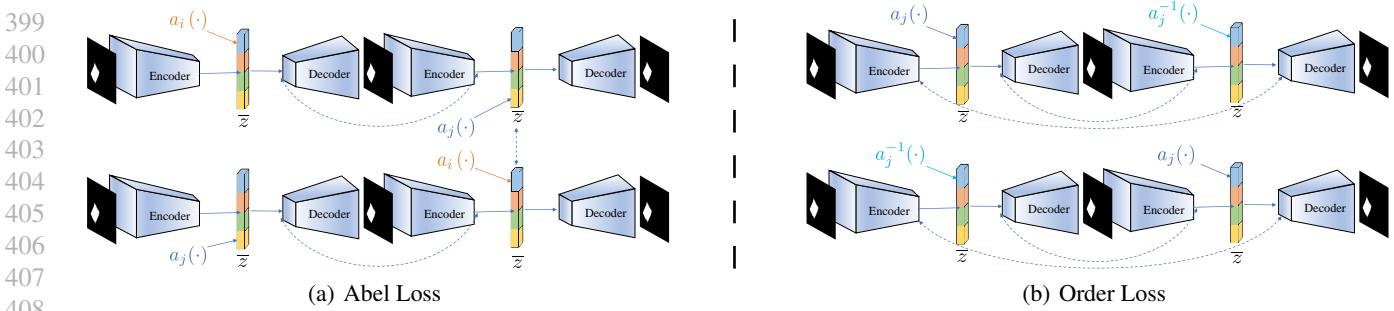
## 385 E. Details of Implementation

### 386 E.1. Abel Loss Details

388 As mentioned in the main paper, the Abel Loss of the groupifiable VAE-based model is as follows:

$$389 \quad 390 \quad 391 \quad 392 \quad \mathcal{L}_a = \sum_{x \in X} \sum_{(i,j) \in C} \|\varphi_i(\varphi_j(x)) - \varphi_j(\varphi_i(x))\|, \quad (40)$$

393 where  $\varphi_i(\varphi_j(x))$  represents the upper path of Figure 3 (a), and  $\varphi_j(\varphi_i(x))$  represents the lower path of Figure 3 (a). For  
394 better optimization, we do not constrain such consistency on the reconstructed images but on their representations instead  
395 (straight dotted line in Figure 3 (a)). Besides, we constrain the consistency between the representations of intermediate  
396 images (curved dotted lines in Figure 3 (a)).  $a_i$  denotes the group action of  $e_i$  on  $Z$ , i.e.,  $a_i(\bar{z}) = z + e_i$ . We implement  $a_i$   
397 by adding the action scale  $\tau$  on the  $i$ -th dimension of  $z$ , then mapping it to  $\bar{z}$  by the sine and cosine functions, so as to  $a_i$ .



400 Figure 3. Overview of the Isomorphism Loss. The Abel Loss constrains the exchangeability of Permutation Group  $\Phi$ . The Order Loss  
401 constrains the cyclicity of Permutation Group  $\Phi$ . Dotted line in the figure represent reconstruction loss.  $a_i, a_j, a_i^{-1}$  denotes group actions.

### 402 E.2. Order Loss Details

403 As mentioned in the main paper, the Order Loss of the groupified VAEs is as follows:

$$404 \quad 405 \quad 406 \quad \mathcal{L}_o = \sum_{x \in X} \sum_{i \in I} (\|\varphi_i(\varphi_i^{-1}(x)) - x\| + \|\varphi_i^{-1}(\varphi_i(x)) - x\|), \quad (41)$$

407 where  $\varphi_i(\varphi_i^{-1}(x))$  represents the lower path of Figure 3 (b), and  $\varphi_i^{-1}(\varphi_i(x))$  represents the upper path of Figure 3 (b).  
408 Similar to the Abel Loss, we do not constrain such consistency on the reconstructed images for better optimization, but on  
409 their representations instead (long curved dotted line in Figure 3 (b)). Besides, we constrain the consistency between the  
410 representations of intermediate images (short curved dotted lines in Figure 3 (b)).  $a_j^{-1}$  denotes the group action of  $e_j^{-1}$  on  $Z$ ,  
411 i.e.,  $a_j^{-1}(\bar{z}) = z + (n-1)e_j$ . We implement  $a_j^{-1}$  by adding the action scale  $(N-1)\tau$  on the  $j$ -th dimension of  $z$ , then  
412 mapping it to  $\bar{z}$  by the sine and cosine functions (we set  $n = N$ ).

### 413 E.3. The set of factors $I$

414 For the given VAE-based models, there are two ways to get set  $I$ . (i) Assign some dimensions to  $I$ , e.g., set  $I = \{0, 1, 2, 3, 4\}$ ,  
415 which is the first-5 of the dimensions, resulting in meaningful dimensions controllable. (ii) Dimensional KL divergence  
416 indicates the meaningful dimensions. We set  $I = \{i | KL_i \geq T\}$ , where  $T$  is a hyperparameter, which is set to 30 in our  
417 experiments.

## 440 F. Details of Experiments

## 441 F1. Dataset details

443 In all the experiments, we resize the images to 64x64 resolution. We then introduce all the datasets used in our paper in  
444 details.

**dSprites** ([Higgins et al., 2016](#)): dSprites contains 737,280 binary 2D shapes (heart, oval and square) images with 4 ground truth factors: shape (3 values), scale (6 values), orientation (40 values), x-position (32 values), y-position (32 values). Then we introduce the variants of dSprites created by [Locatello et al. \(2019\)](#).

449 **Color-dSprites**, the shapes are colored with a random color.

**Noisy-dSprites**, we consider white-colored shapes on a noisy background.

**Shapes3D (Kim & Mnih, 2018):** Shapes3D contains 480,000 3D shapes images with 6 ground truth factors: shape (4 values), scale (8 values), orientation (15 values), floor color (10 values), wall color (15 values), object color (10 values).

**Cars3D (Reed et al., 2015)**: This dataset consists of 183 car CAD models, each rendered from 24 azimuth directions and 4 elevations.

Table 1. Architecture of the encoder and decoder of VAEs. For original VAE, the dimension of input of the decoder is 10. For Groupified VAE, the dimension is 20.

Encoder	Decoder
Input: $64 \times 64 \times$ number of channels	Input: 10 or 20
$4 \times 4$ conv, 32 ReLU, stride 2	FC, 256 ReLU
$4 \times 4$ conv, 32 ReLU, stride 2	FC, 256 ReLU
$4 \times 4$ conv, 32 ReLU, stride 2	FC, $4 \times 4 \times 32$ ReLU
$4 \times 4$ conv, 32 ReLU, stride 2	$4 \times 4$ deconv, 32 ReLU, stride 2
FC 256 ReLU	$4 \times 4$ deconv, 32 ReLU, stride 2
FC 256 ReLU	$4 \times 4$ deconv, 32 ReLU, stride 2
FC $2 \times 10$	$4 \times 4$ deconv, number of channels, stride 2

## 472 F.2. Architecture for encoder and decoder

We follow Locatello et al. (2019) to use the same architecture of VAEs in all of the experiments as follows: the activation function used is ReLU except for the last layer of decoder, which is Sigmoid. As shown in Table 1. For the details of the Discriminator in FactorVAE, please refers to Table 3 (a) and (c).

Table 2. Hyperparameters and random seeds for every model.

Model	Parameter	Value
$\beta$ -VAE	$\beta$	[10; 20; 30;]
AnnealedVAE	C	[10; 20; 30;]
	start	[3e4; 4e4; 5e4;]
	end	[2e4; 3e4; 4e4;]
FactorVAE	$\gamma$	[5; 10; 15]
$\beta$ -TCVAE	$\beta$	[6; 9; 12]
	random seed	[1; 2; 3; 4; 5; 6; 7; 8; 9;]
	group	[True; False]

### 491 F.3. Experiment settings

We run different hyperparameters and random seeds for every VAE-based models implemented by Pytorch ([Paszke et al., 2017](#)). As shown in Table 2, for  $\beta$ -VAE, we assign 3 choices for  $\beta$  and 10 random seeds for both the original and groupified

Table 3. Shared hyperparameters in all experiments.

(a) Optimizer for Discriminator		(b) General hyperparameters for VAE		(c) Architecture of Discriminator
Parameter	Values	Parameter	Values	Discriminator
Batch size	64	Batch size	64	FC, 1000 leaky ReLU
Optimizer	Adam	Latent space dimension	10	FC, 1000 leaky ReLU
Adam: beta1	0.9	Optimizer	Adam	FC, 1000 leaky ReLU
Adam: beta2	0.999	Adam: beta1	0.9	FC, 1000 leaky ReLU
Adam: epsilon	1.00e-08	Adam: beta2	0.999	FC, 1000 leaky ReLU
Adam: learning rate	0.0001	Adam: epsilon	1.00e-08	FC, 1000 leaky ReLU
		Adam: learning rate	0.0001	FC, 2
		Decoder type	Bernoulli	

VAEs:  $3 \times 10 \times 2 = 60$  settings for each dataset. Similarly, we also assign 60 settings for FactorVAE and  $\beta$ -TCVAE. For AnnealVAE, we assign 3 choices for  $C_{max}$  and 3 choices for the start and end pair, also assign 10 random seeds. In summary, for all 5 datasets, we run  $((3 \times 10 \times 2) \times 3) + 3 \times 3 \times 10 \times 2 \times 5 = 1800$  models. For other hyperparameters, please refer to Table 3 (b).

## G. More Qualitative Results and Score Distribution

### G.1. Qualitative Results

The performance of original and groupified VAEs on all 5 datasets is shown in Table 4. Our method outperforms the original one on most of the datasets in terms of nearly all the metrics.

*Table 4.* Performance (mean  $\pm$  variance) on different datasets and different models with different metrics. We evaluate  $\beta$ -VAE, AnnealVAE, FactorVAE, and  $\beta$ -TCVAE on dSprites, Cars3d, Shapes3d, Noisy-dSprites, and Color-dSprites for 1800 settings. These settings include different random seeds and hyperparameters.

dSprites	DCI		BetaVAE		MIG		FactorVAE	
	Original	Groupified	Original	Groupified	Original	Groupified	Original	Groupified
$\beta$ -VAE	0.23 $\pm$ 0.10	<b>0.46 <math>\pm</math> 0.085</b>	0.75 $\pm$ 0.083	<b>0.86 <math>\pm</math> 0.051</b>	0.14 $\pm$ 0.097	<b>0.37 <math>\pm</math> 0.089</b>	0.51 $\pm$ 0.098	<b>0.63 <math>\pm</math> 0.089</b>
AnnealVAE	0.28 $\pm$ 0.10	<b>0.39 <math>\pm</math> 0.056</b>	0.84 $\pm$ 0.050	<b>0.87 <math>\pm</math> 0.0067</b>	0.23 $\pm$ 0.10	<b>0.34 <math>\pm</math> 0.061</b>	<b>0.70 <math>\pm</math> 0.094</b>	0.68 $\pm$ <b>0.058</b>
FactorVAE	0.38 $\pm$ 0.10	<b>0.41 <math>\pm</math> 0.074</b>	<b>0.89 <math>\pm</math> 0.040</b>	<b>0.89 <math>\pm</math> 0.020</b>	0.27 $\pm$ 0.092	<b>0.31 <math>\pm</math> 0.061</b>	0.74 $\pm$ 0.068	<b>0.75 <math>\pm</math> 0.075</b>
$\beta$ -TCVAE	0.35 $\pm$ 0.065	<b>0.36 <math>\pm</math> 0.11</b>	0.86 $\pm$ 0.026	<b>0.861 <math>\pm</math> 0.038</b>	0.031 $\pm$ 0.17	<b>0.060 <math>\pm</math> 0.24</b>	0.68 $\pm$ 0.098	<b>0.70 <math>\pm</math> 0.098</b>
Cars3d	DCI		BetaVAE		MIG		FactorVAE	
	Original	Groupified	Original	Groupified	Original	Groupified	Original	Groupified
$\beta$ -VAE	0.18 $\pm$ 0.059	<b>0.24 <math>\pm</math> 0.041</b>	0.99 $\pm$ 1.6e $-3$	<b>1.0 <math>\pm</math> 0.0</b>	0.071 $\pm$ 0.032	<b>0.11 <math>\pm</math> 0.032</b>	0.81 $\pm$ 0.066	<b>0.93 <math>\pm</math> 0.034</b>
AnnealVAE	0.22 $\pm$ 0.046	<b>0.25 <math>\pm</math> 0.046</b>	<b>0.99 <math>\pm</math> 4e <math>-4</math></b>	<b>0.99 <math>\pm</math> 1.5e <math>-4</math></b>	0.074 $\pm$ 0.016	<b>0.10 <math>\pm</math> 0.014</b>	0.82 $\pm$ 0.062	<b>0.87 <math>\pm</math> 0.028</b>
FactorVAE	0.21 $\pm$ 0.054	<b>0.25 <math>\pm</math> 0.040</b>	0.99 $\pm$ 1e $-4$	<b>1.0 <math>\pm</math> 0.0</b>	0.098 $\pm$ 0.027	<b>0.11 <math>\pm</math> 0.033</b>	0.90 $\pm$ 0.039	<b>0.93 <math>\pm</math> 0.034</b>
$\beta$ -TCVAE	0.24 $\pm$ 0.049	<b>0.26 <math>\pm</math> 0.046</b>	<b>1.0 <math>\pm</math> 0.0</b>	<b>1.0 <math>\pm</math> 0.0</b>	0.10 $\pm$ 0.021	<b>0.11 <math>\pm</math> 0.033</b>	0.88 $\pm$ 0.040	<b>0.93 <math>\pm</math> 0.034</b>
Noisy dSprites	DCI		betaVAE		MIG		FactorVAE	
	Original	Groupified	Original	Groupified	Original	Groupified	Original	Groupified
BetaVAE	0.056 $\pm$ 0.018	<b>0.087 <math>\pm</math> 0.051</b>	0.624 $\pm$ 0.090	<b>0.647 <math>\pm</math> 0.055</b>	0.030 $\pm$ 0.022	<b>0.065 <math>\pm</math> 0.055</b>	0.355 $\pm$ 0.093	<b>0.407 <math>\pm</math> 0.071</b>
Anneal VAE	0.053 $\pm$ 0.013	<b>0.060 <math>\pm</math> 0.022</b>	0.631 $\pm$ 0.036	<b>0.644 <math>\pm</math> 0.031</b>	0.035 $\pm$ 0.027	<b>0.047 <math>\pm</math> 0.032</b>	0.434 $\pm$ 0.080	<b>0.481 <math>\pm</math> 0.087</b>
FactorVAE	<b>0.114 <math>\pm</math> 0.062</b>	0.099 $\pm$ <b>0.057</b>	0.682 $\pm$ 0.081	<b>0.684 <math>\pm</math> 0.070</b>	<b>0.077 <math>\pm</math> 0.046</b>	0.066 $\pm$ <b>0.046</b>	0.437 $\pm$ 0.098	<b>0.468 <math>\pm</math> 0.098</b>
BetaTCVAE	0.081 $\pm$ 0.036	<b>0.111 <math>\pm</math> 0.053</b>	0.605 $\pm$ 0.053	<b>0.635 <math>\pm</math> 0.050</b>	0.040 $\pm$ 0.030	<b>0.068 <math>\pm</math> 0.042</b>	0.353 $\pm$ 0.091	<b>0.431 <math>\pm</math> 0.097</b>
Shapes3d	DCI		BetaVAE		MIG		FactorVAE	
	Original	Groupified	Original	Groupified	Original	Groupified	Original	Groupified
$\beta$ -VAE	0.44 $\pm$ 0.176	<b>0.56 <math>\pm</math> 0.10</b>	0.91 $\pm$ 0.072	<b>0.90 <math>\pm</math> 0.045</b>	0.28 $\pm$ 0.18	<b>0.42 <math>\pm</math> 0.15</b>	<b>0.82 <math>\pm</math> 0.098</b>	<b>0.82 <math>\pm</math> 0.043</b>
AnnealVAE	0.52 $\pm$ 0.051	<b>0.60 <math>\pm</math> 0.078</b>	0.82 $\pm$ <b>0.076</b>	<b>0.89 <math>\pm</math> 0.086</b>	0.48 $\pm$ 0.047	<b>0.50 <math>\pm</math> 0.052</b>	0.75 $\pm$ 0.074	<b>0.83 <math>\pm</math> 0.066</b>
FactorVAE	0.47 $\pm$ 0.10	<b>0.49 <math>\pm</math> 0.065</b>	<b>0.86 <math>\pm</math> 0.055</b>	0.80 $\pm$ 0.075	0.33 $\pm$ 0.13	<b>0.43 <math>\pm</math> 0.11</b>	<b>0.81 <math>\pm</math> 0.056</b>	0.79 $\pm$ 0.066
$\beta$ -TCVAE	0.66 $\pm$ 0.10	<b>0.72 <math>\pm</math> 0.061</b>	<b>0.97 <math>\pm</math> 0.039</b>	0.96 $\pm$ 0.042	0.40 $\pm$ 0.18	<b>0.47 <math>\pm</math> 0.090</b>	0.89 $\pm$ 0.064	<b>0.90 <math>\pm</math> 0.046</b>
Color dSprites	DCI		betaVAE		MIG		FactorVAE	
	Original	Groupified	Original	Groupified	Original	Groupified	Original	Groupified
BetaVAE	0.174 $\pm$ 0.097	<b>0.328 <math>\pm</math> 0.130</b>	0.798 $\pm$ 0.094	<b>0.844 <math>\pm</math> 0.050</b>	0.103 $\pm$ 0.058	<b>0.243 <math>\pm</math> 0.118</b>	0.591 $\pm$ 0.148	<b>0.648 <math>\pm</math> 0.092</b>
Anneal VAE	0.268 $\pm$ 0.103	<b>0.337 <math>\pm</math> 0.114</b>	0.843 $\pm$ 0.038	<b>0.856 <math>\pm</math> 0.031</b>	0.219 $\pm$ 0.084	<b>0.252 <math>\pm</math> 0.104</b>	<b>0.718 <math>\pm</math> 0.065</b>	0.692 $\pm$ 0.094
FactorVAE	0.294 $\pm$ 0.101	<b>0.322 <math>\pm</math> 0.104</b>	0.861 $\pm$ 0.038	<b>0.862 <math>\pm</math> 0.029</b>	0.203 $\pm$ 0.080	<b>0.236 <math>\pm</math> 0.091</b>	<b>0.739 <math>\pm</math> 0.068</b>	0.730 $\pm$ 0.080
BetaTCVAE	0.338 $\pm$ 0.052	<b>0.395 <math>\pm</math> 0.082</b>	0.876 $\pm$ 0.024	<b>0.881 <math>\pm</math> 0.031</b>	0.169 $\pm$ 0.040	<b>0.269 <math>\pm</math> 0.090</b>	0.711 $\pm$ 0.086	<b>0.786 <math>\pm</math> 0.050</b>

### G.2. Score Distribution

The detailed distribution of the performance is shown in this section (demonstrated by the Violin Plot (Hintze & Nelson, 1998)). The performance distributions on dSprites, Car3d, Noisy dSprites, Color dSprites and Shapes3d are shown in Figure 4, Figure 5, Figure 6, Figure 7 and Figure 8 respectively.

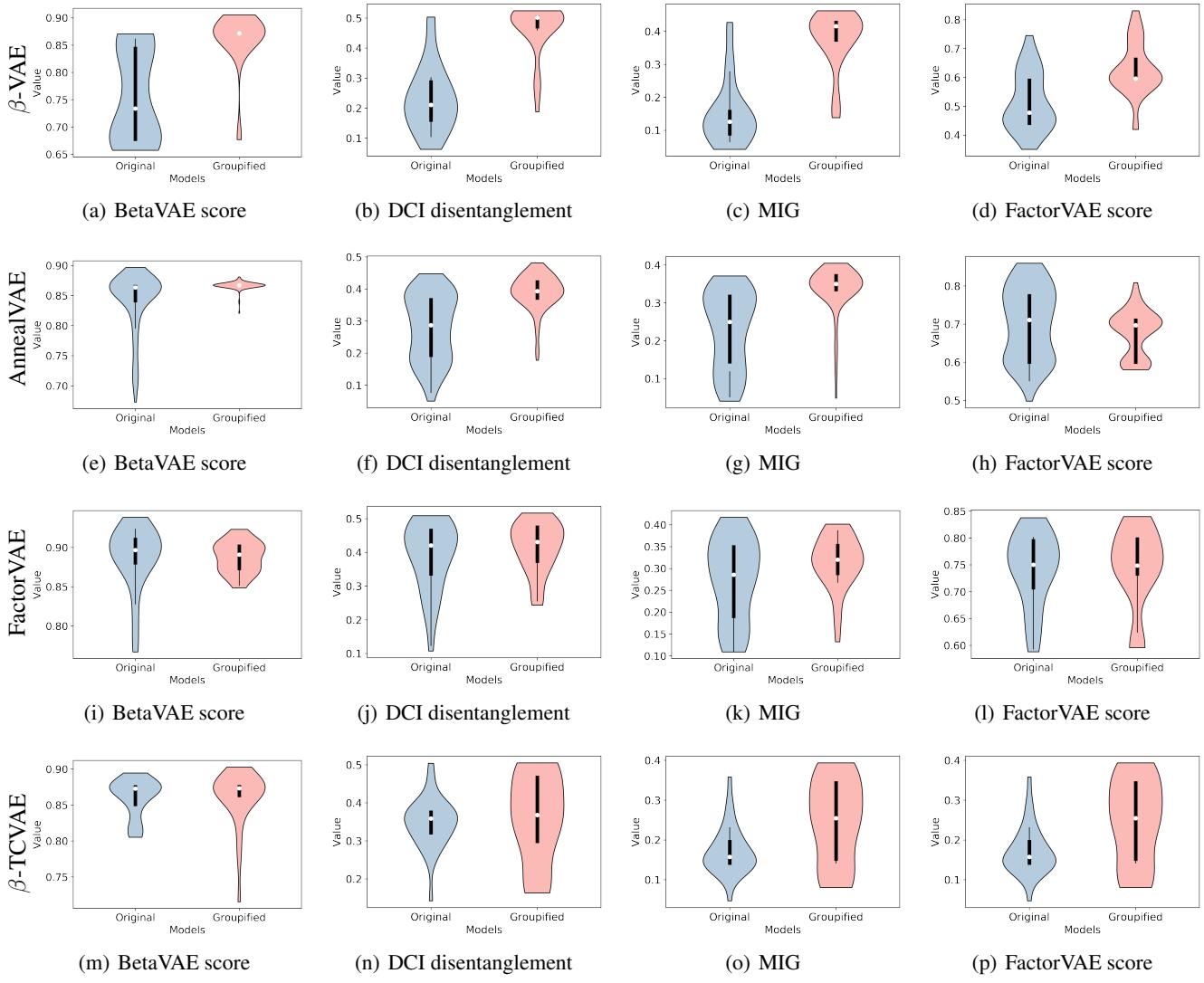


Figure 4. Performance distribution on dSprites. Variance is due to different hyperparameters and random seeds. We take four metrics into consideration: BetaVAE score, DCI disentanglement, MIG and FactorVAE score. We observe that groupified models outperform the original ones.

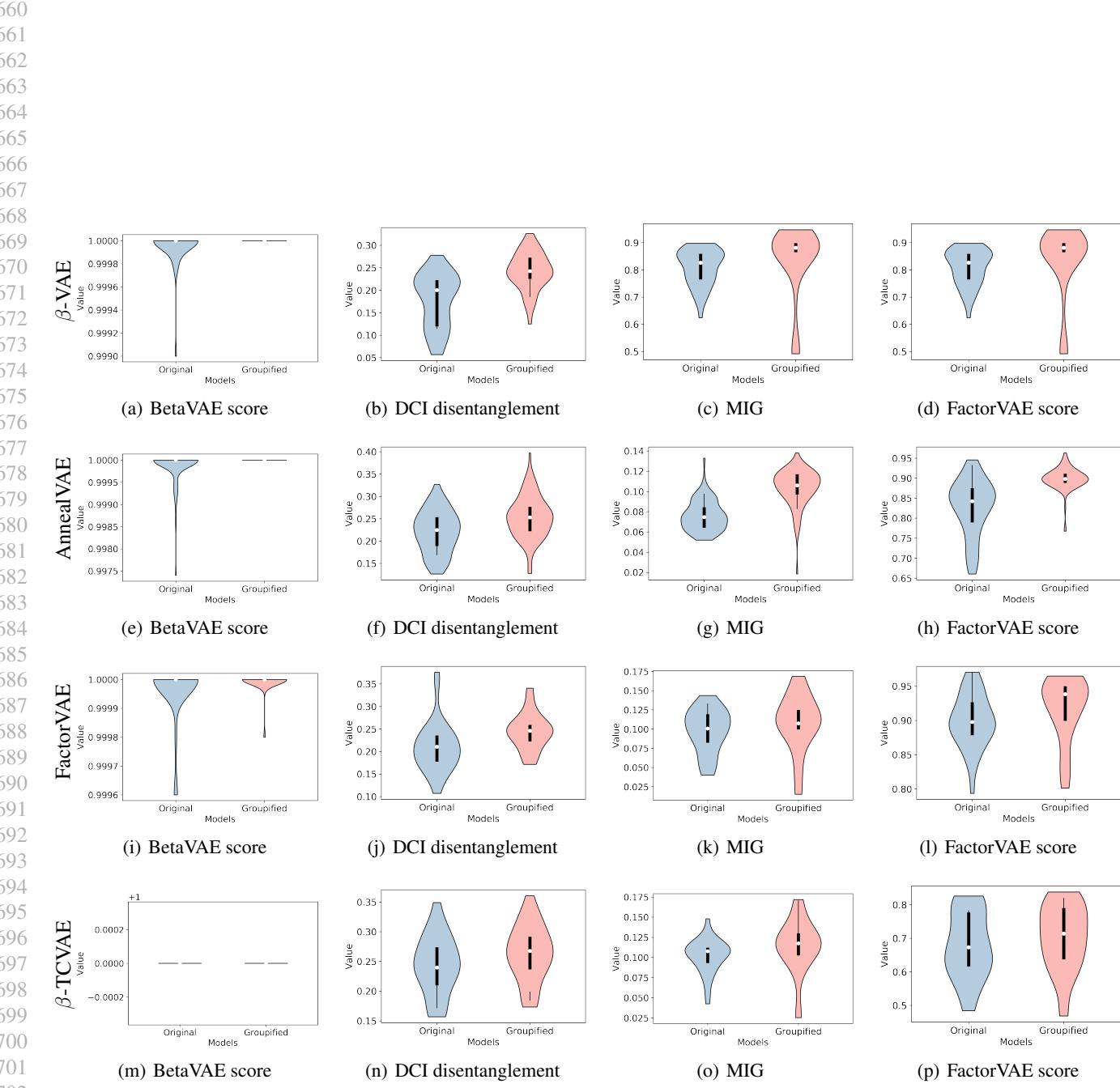


Figure 5. Performance distribution on Cars3d. Variance is due to different hyperparameters and random seeds. We take four metrics into consideration: BetaVAE score, DCI disentanglement, MIG and FactorVAE score. We observe that groupified models outperform the original ones.

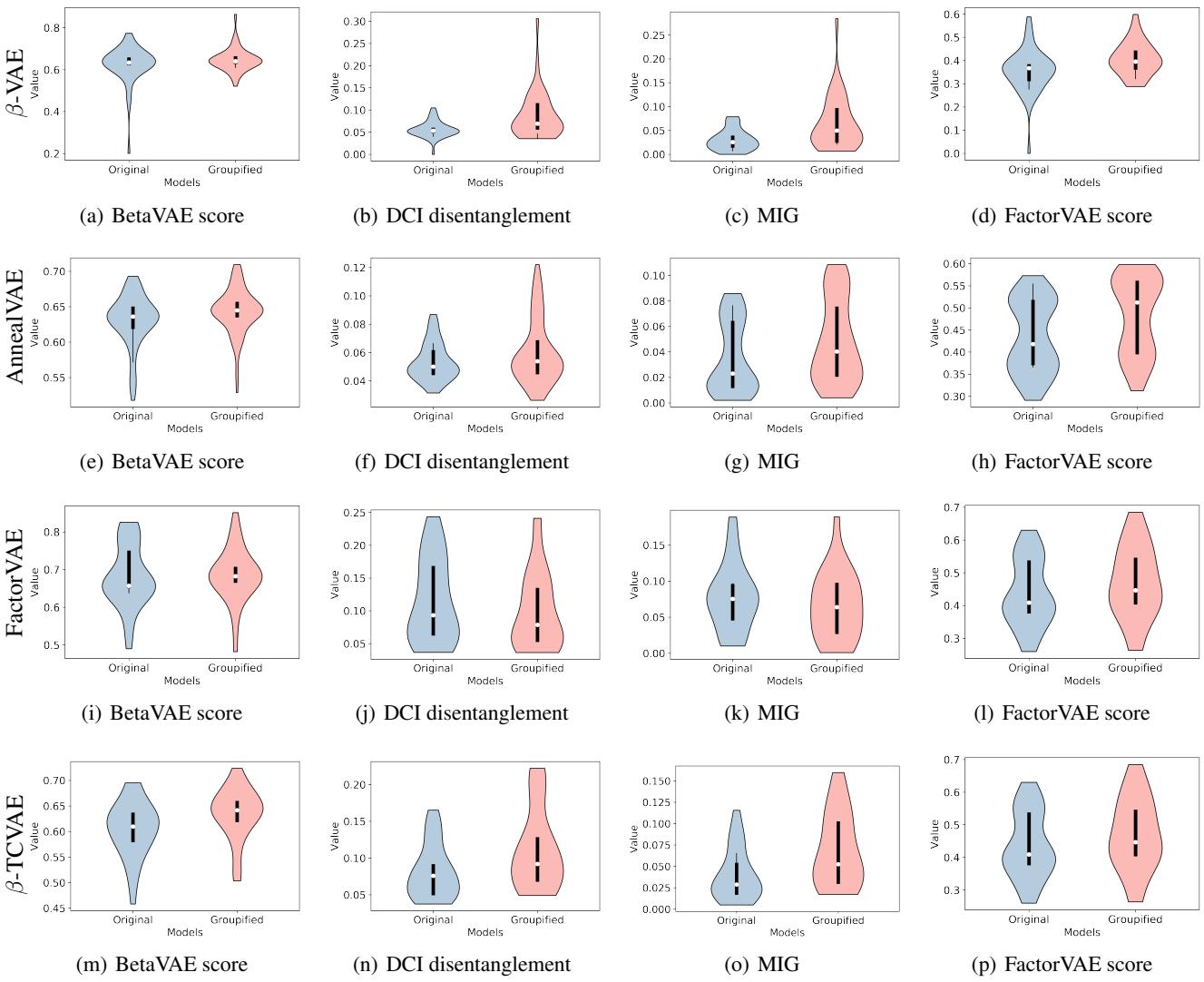


Figure 6. Performance distribution on Noisy dSprites. Variance is due to different hyperparameters and random seeds. We take four metrics into consideration: BetaVAE score, DCI disentanglement, MIG and FactorVAE score. We observe that groupified models outperform the original ones.

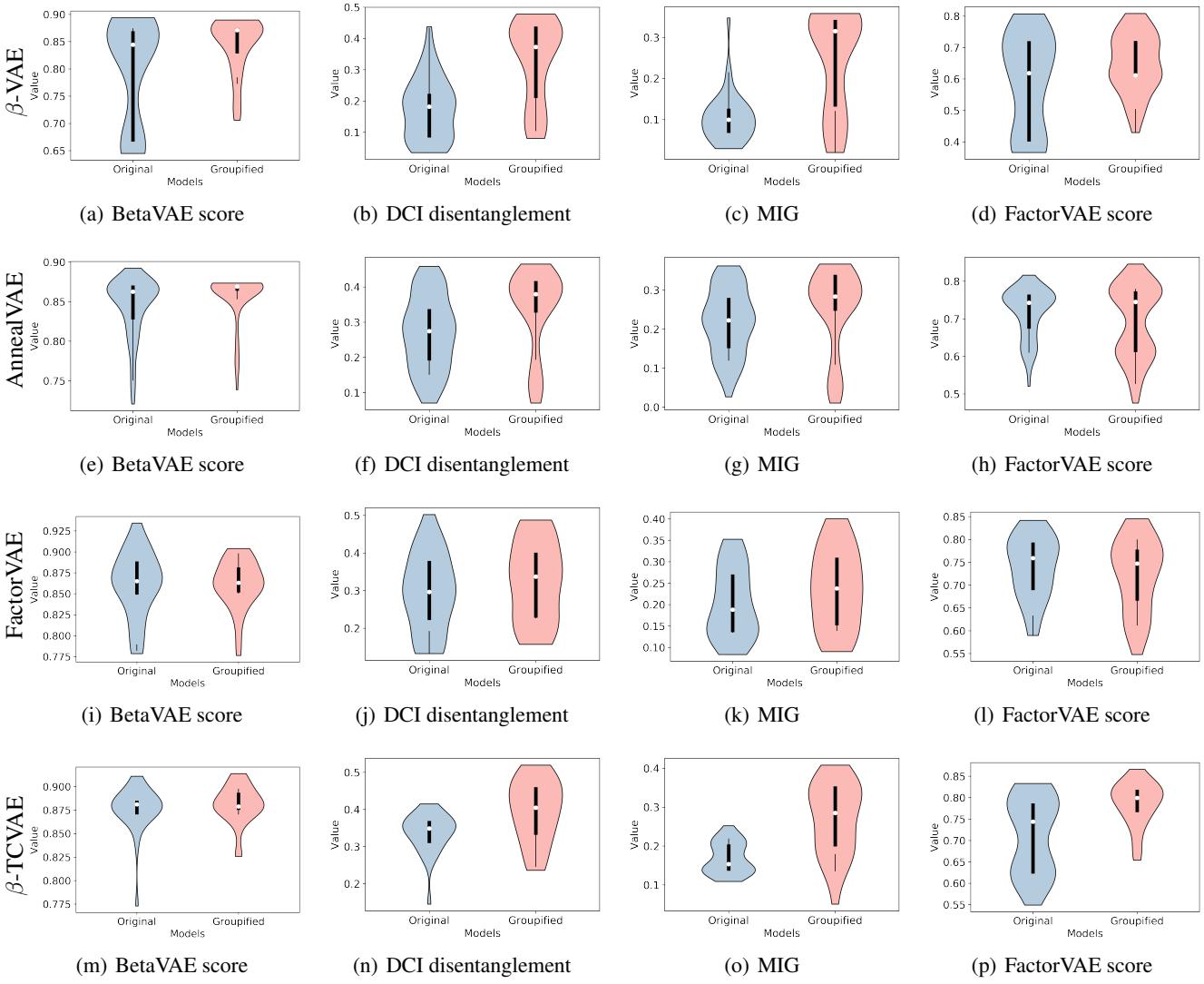


Figure 7. Performance distribution on Color dSprites. Variance is due to different hyperparameters and random seeds. We take four metrics into consideration: BetaVAE score, DCI disentanglement, MIG and FactorVAE score. We observe that groupified models outperform the original ones.

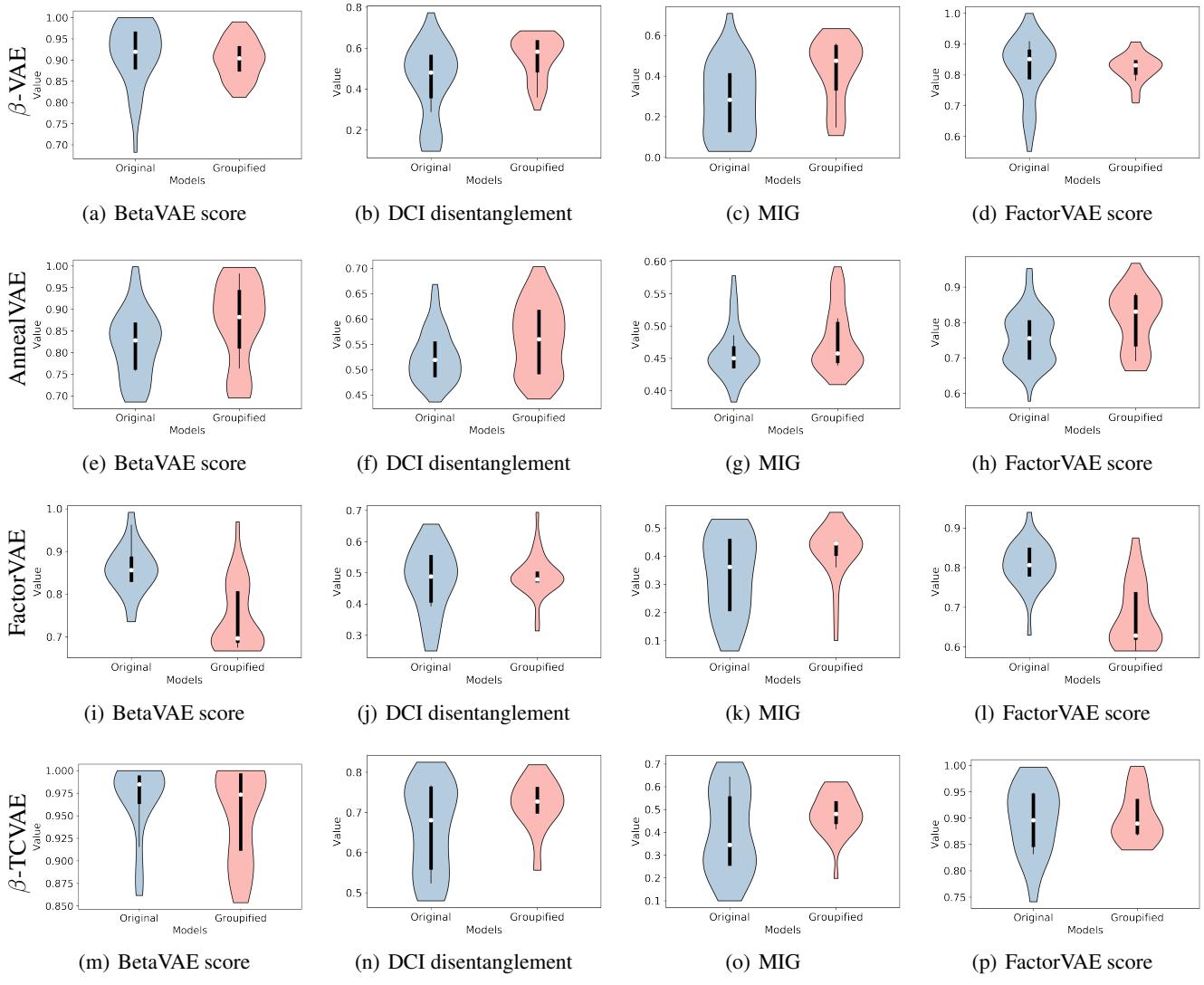
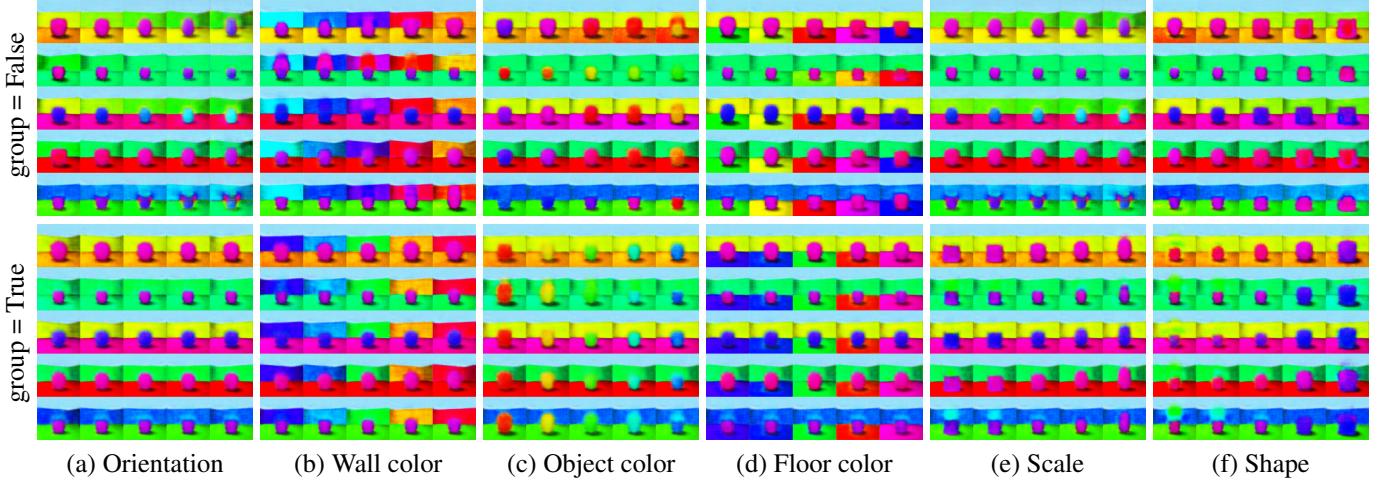


Figure 8. Performance distribution on Shapes3d. Variance is due to different hyperparameters and random seeds. We take four metrics into consideration: BetaVAE score, DCI disentanglement, MIG and FactorVAE score. We observe that groupified models outperform the original ones.

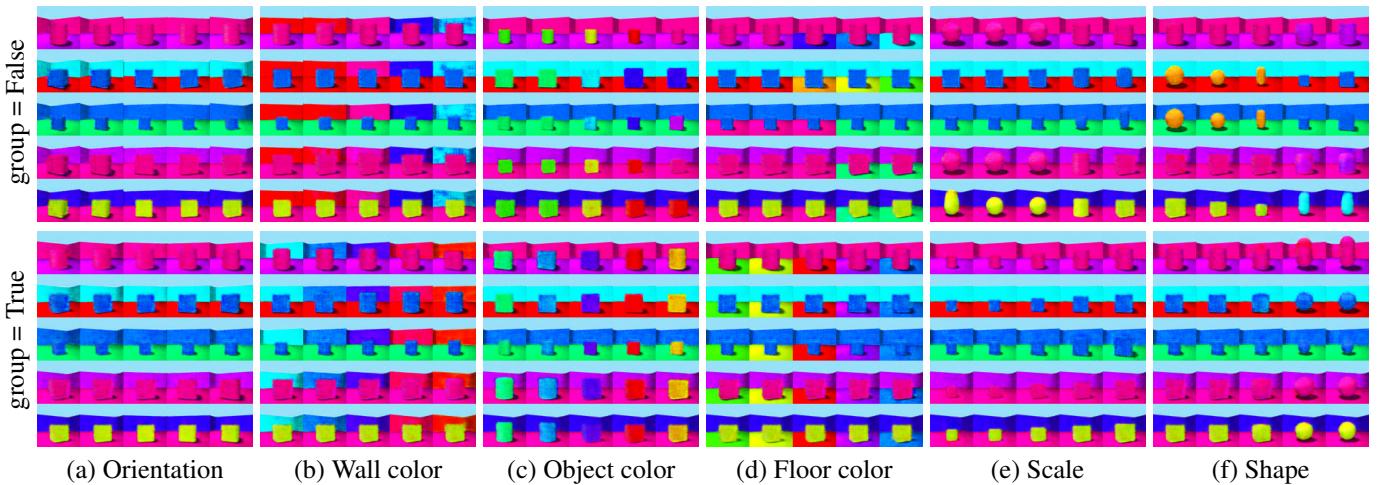
## 880 H. More Qualitative Results

881  
882 We evaluate our methods qualitatively on two typical datasets: Cars3d and Shapes3d. We visualize the traversal results of  
883 orginal and groupified FactorVAE and  $\beta$ -TCVAE. For every factor, we traversal from 5 random representation. As shown in  
884 Figure 9 and Figure 10, the traversal results of groupified FactorVAE and  $\beta$ -TCVAE shows that these models learn less  
885 entangled representation on Shapes3D (e.g., Orientation of FactorVAE and Scale and Shape of  $\beta$ -TCVAE).

886 Similarly, as shown in Figure 11 and Figure 12, groupified FactorVAE and  $\beta$ -TCVAE archive better disentanglement ability  
887 on Car3d (e.g., Rotation of FactorVAE and Yaw of  $\beta$ -TCVAE).



905 Figure 9. Learned latent variables using original and groupified FactorVAE on Shapes3d dataset. Traversal range is (-2, 2).  
906



924 Figure 10. Learned latent variables using original and groupified  $\beta$ -TCVAE on Shapes3d dataset. Traversal range is (-2, 2).  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934

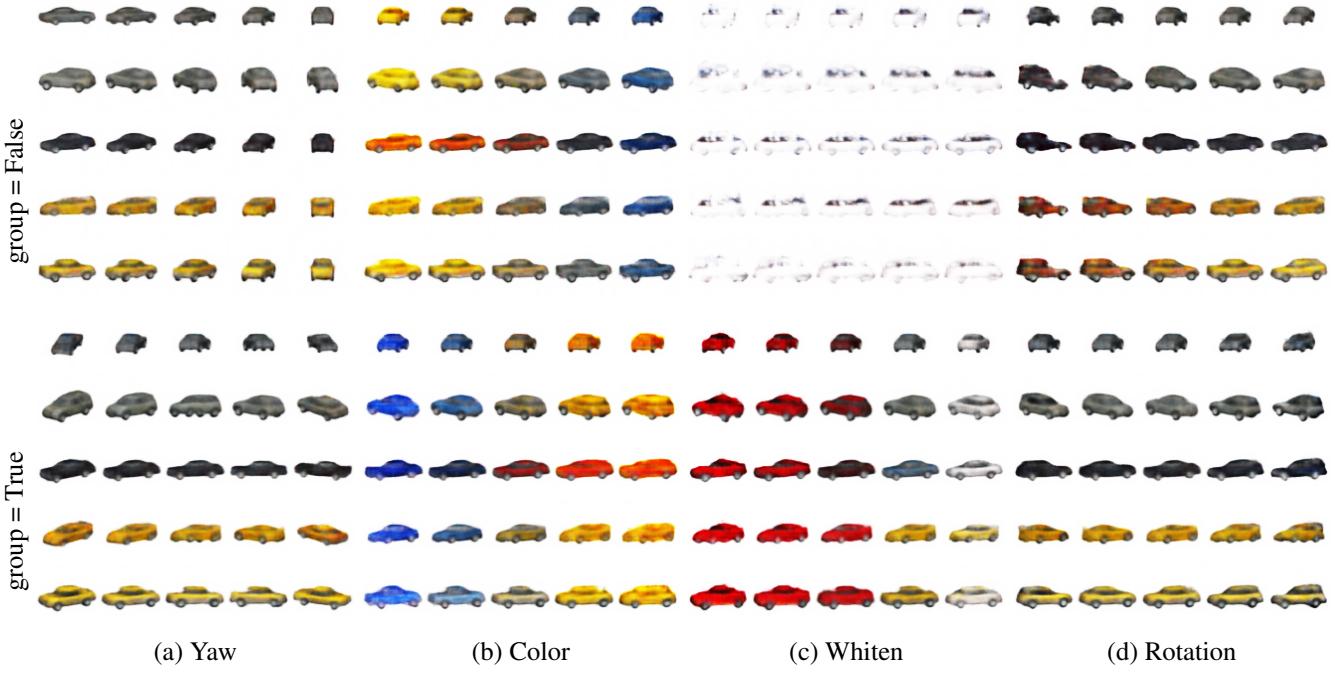


Figure 11. Learned latent variables using original and groupified FactorVAE on Car3d dataset. Traversal range is (-2, 2).

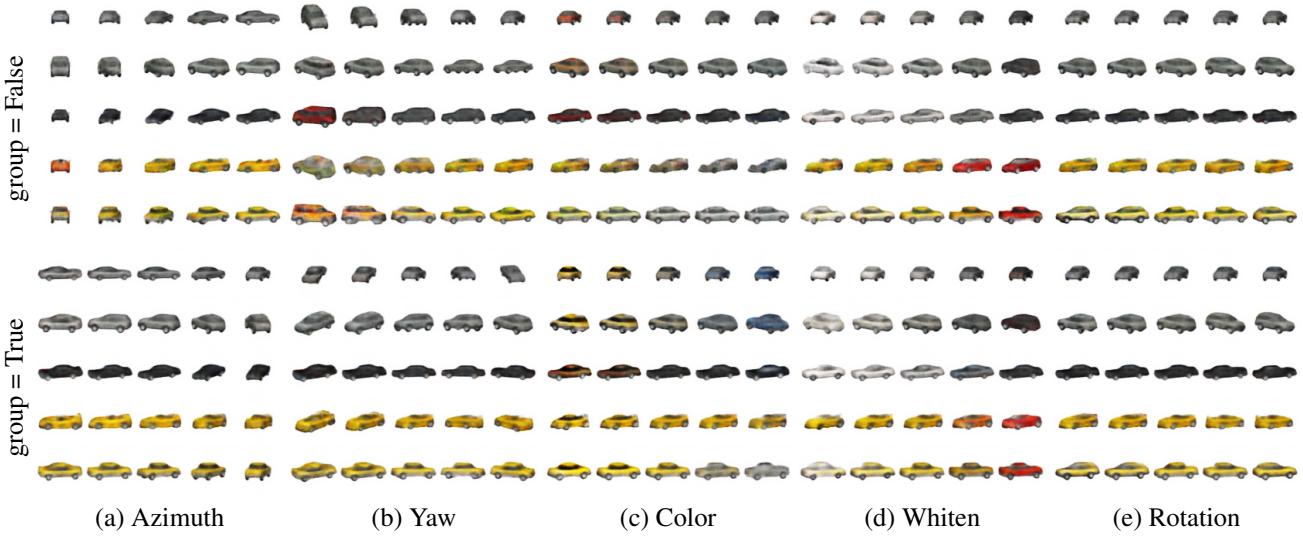
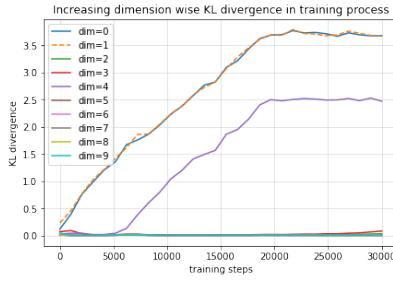


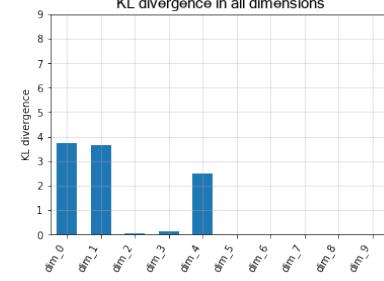
Figure 12. Learned latent variables using original and groupified  $\beta$ -TCVAE on Car3d dataset. Traversal range is (-2, 2).

## I. More Meaningful Dimension Visualizations

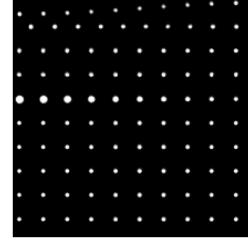
To illustrate that controllable dimensions in AnnealVAE is not an exception, we provide more visualizations of groupified AnnealVAEs. Groupified AnnealVAEs of different hyperparameters and random seeds are shown in Figure 13 to 21, suggesting that the dimensions of them are controllable.



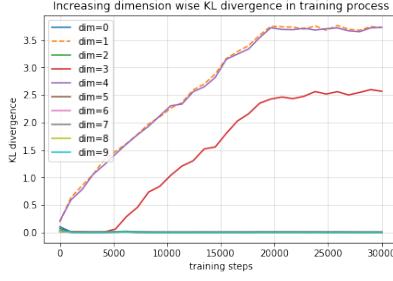
(a) Dimension wised KL in training



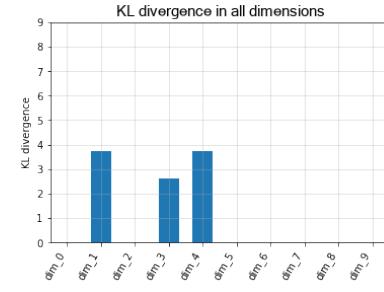
(b) Dimension wise KL in 30000 steps



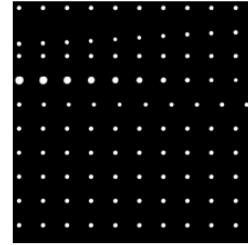
(c) Images traversal of dimensions



(d) Dimension wised KL in training

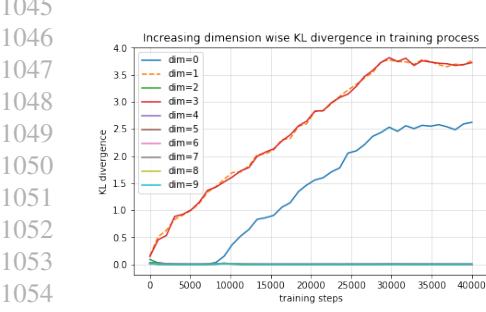


(e) Dimension wise KL in 30000 steps

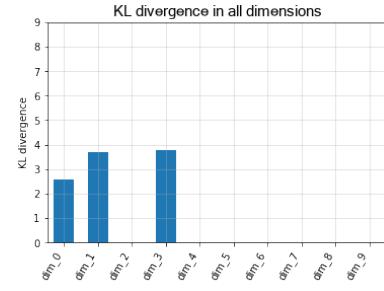


(f) Images traversal of dimensions

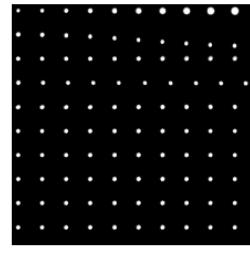
**Figure 13.** Meaningful dimensions visualization for  $C_{max} = 10$ , end = 30000 (different random seeds). The KL divergences of target dimension (0-4 dimension) increase one by one during training (a). The KL divergences in different dimensions are different amounts after training (b). As the the image traversal results (c) shows, the meaningful dimensions are learned in 0-4 dims. So as to (d), (e) and (f).



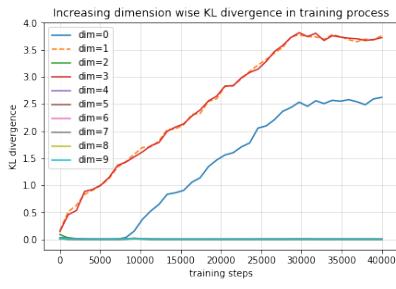
(a) Dimension wised KL in training



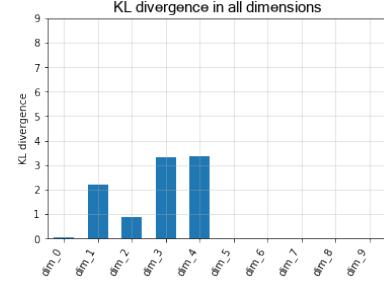
(b) Dimension wise KL in 40000 steps



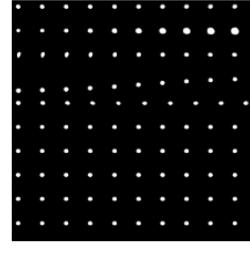
(c) Images traversal of dimensions



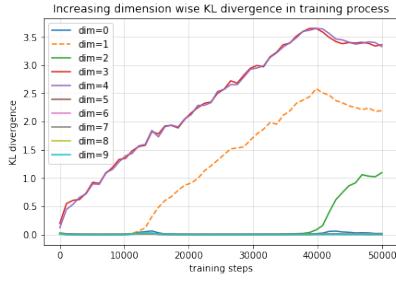
(d) Dimension wised KL in training



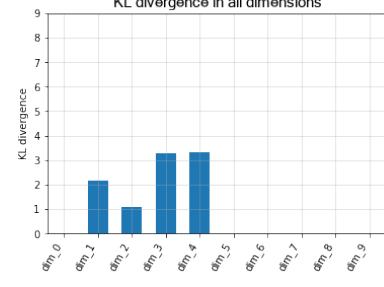
(e) Dimension wise KL in 40000 steps



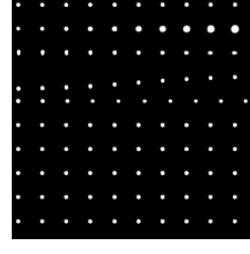
(f) Images traversal of dimensions



(a) Dimension wised KL in training



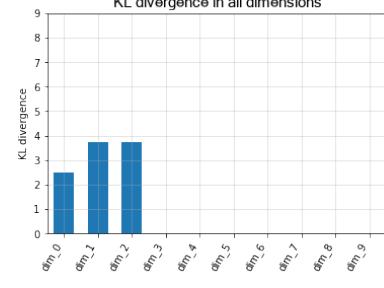
(b) Dimension wise KL in 50000 steps



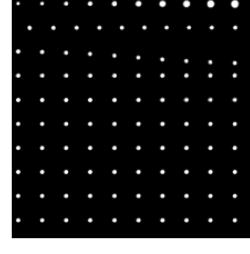
(c) Images traversal of dimensions



(d) Dimension wised KL in training

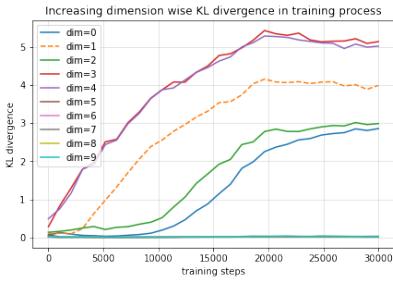


(e) Dimension wise KL in 50000 steps

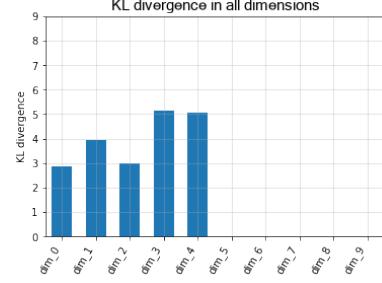


(f) Images traversal of dimensions

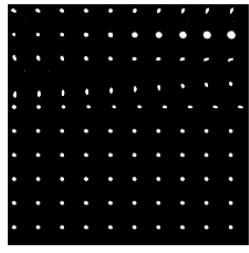
Figure 14. Meaningful dimensions visualization for  $C_{max} = 10$ , end = 40000 (different random seeds). The KL divergences of target dimension (0-4 dimension) increase one by one during training (a). The KL divergences in different dimensions are different amounts after training (b). As the the image traversal results (c) shows, the meaningful dimensions are learned in 0-4 dims. So as to (d), (e) and (f).



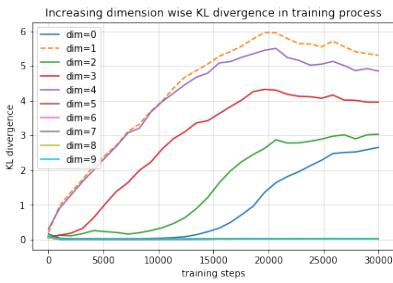
(a) Dimension wised KL in training



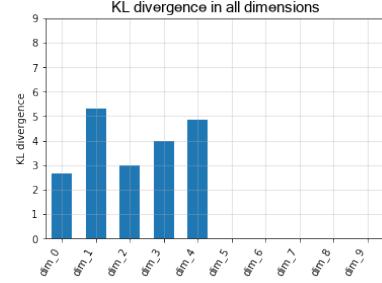
(b) Dimension wise KL in 50000 steps



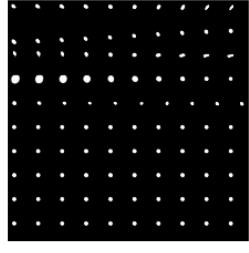
(c) Images traversal of dimensions



(d) Dimension wised KL in training

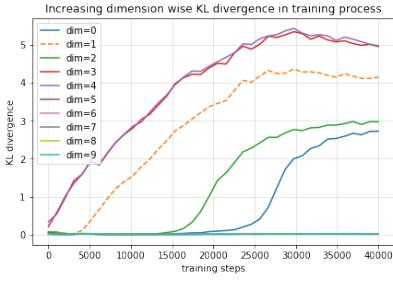


(e) Dimension wise KL in 50000 steps

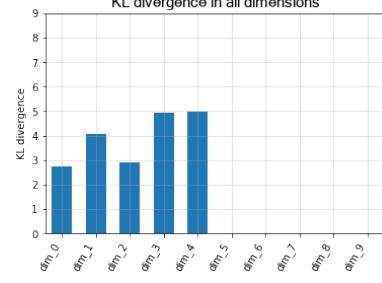


(f) Images traversal of dimensions

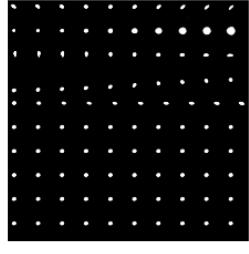
Figure 16. Meaningful dimensions visualization for  $C_{max} = 20$ , end = 30000 (different random seeds). The KL divergences of target dimension (0-4 dimension) increase one by one during training (a). The KL divergences in different dimensions are different amounts after training (b). As the the image traversal results (c) shows, the meaningful dimensions are learned in 0-4 dims. So as to (d), (e) and (f).



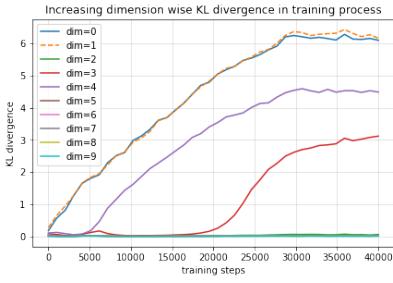
(a) Dimension wised KL in training



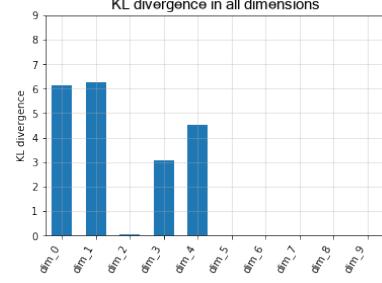
(b) Dimension wise KL in 50000 steps



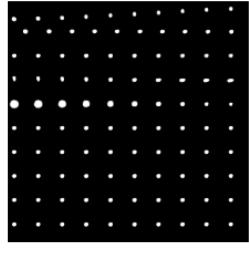
(c) Images traversal of dimensions



(d) Dimension wised KL in training

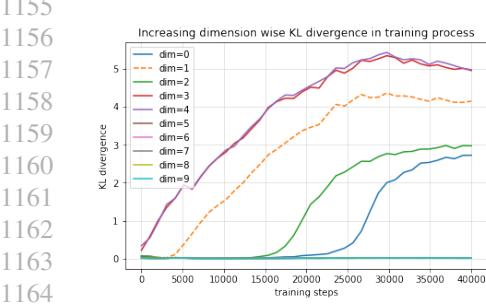


(e) Dimension wise KL in 50000 steps

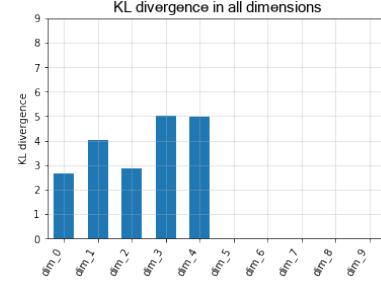


(f) Images traversal of dimensions

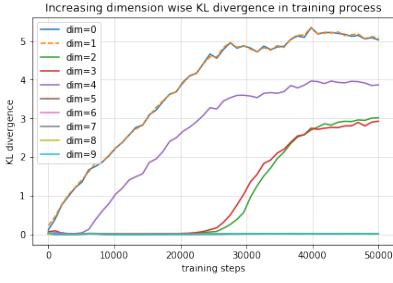
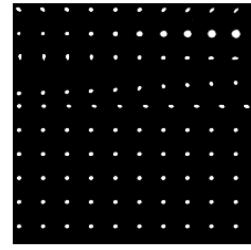
Figure 17. Meaningful dimensions visualization for  $C_{max} = 20$ , end = 40000 (different random seeds). The KL divergences of target dimension (0-4 dimension) increase one by one during training (a). The KL divergences in different dimensions are different amounts after training (b). As the the image traversal results (c) shows, the meaningful dimensions are learned in 0-4 dims. So as to (d), (e) and (f).



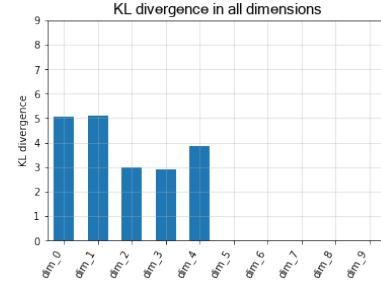
(a) Dimension wised KL in training



(b) Dimension wise KL in 50000 steps



(d) Dimension wised KL in training



(e) Dimension wise KL in 50000 steps

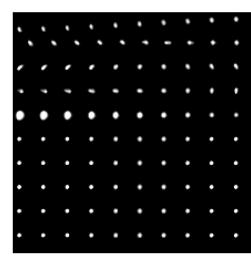
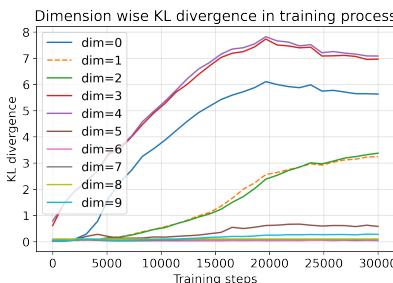
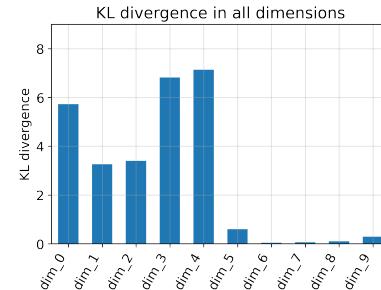


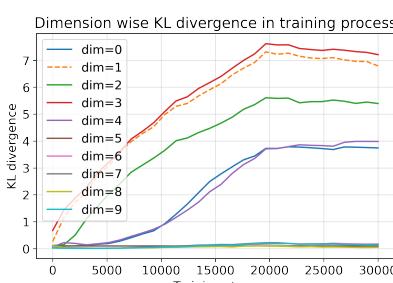
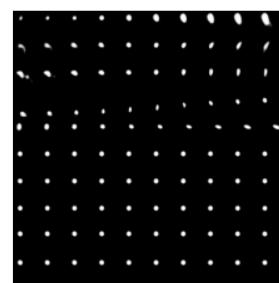
Figure 18. Meaningful dimensions visualization for  $C_{max} = 20$ , end = 50000 (different random seeds). The KL divergences of target dimension (0-4 dimension) increase one by one during training (a). The KL divergences in different dimensions are different amounts after training (b). As the the image traversal results (c) shows, the meaningful dimensions are learned in 0-4 dims. So as to (d), (e) and (f).



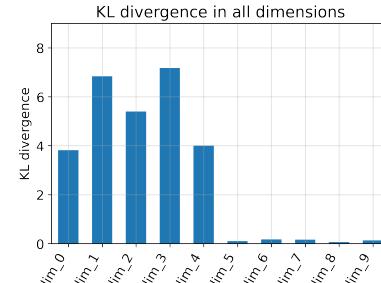
(a) Dimension wised KL in training



(b) Dimension wise KL in 50000 steps



(d) Dimension wised KL in training



(e) Dimension wise KL in 50000 steps

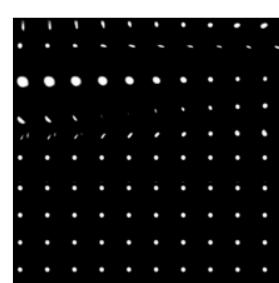
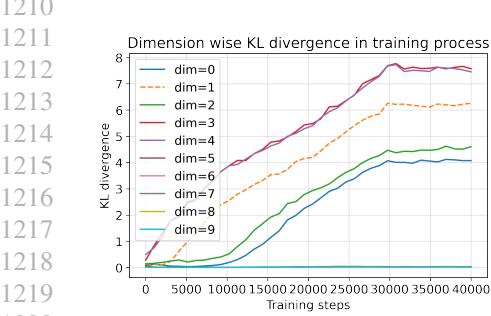
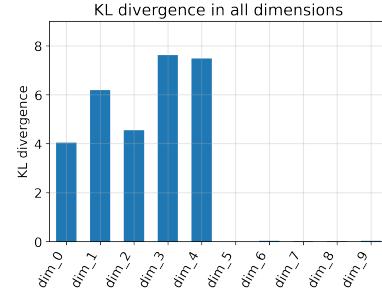


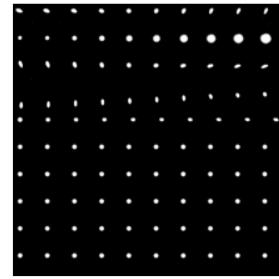
Figure 19. Meaningful dimensions visualization for  $C_{max} = 30$ , end = 30000 (different random seeds). The KL divergences of target dimension (0-4 dimension) increase one by one during training (a). The KL divergences in different dimensions are different amounts after training (b). As the the image traversal results (c) shows, the meaningful dimensions are learned in 0-4 dims. So as to (d), (e) and (f).



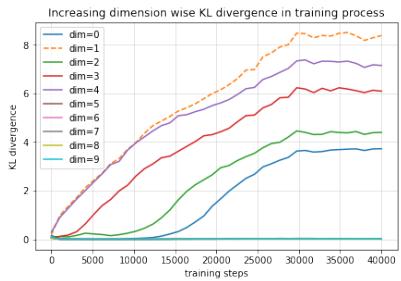
(a) Dimension wised KL in training



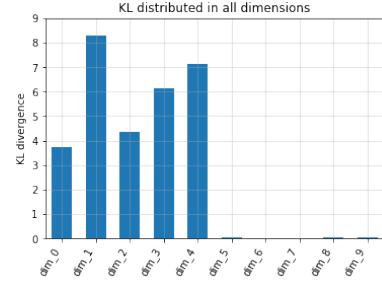
(b) Dimension wise KL in 50000 steps



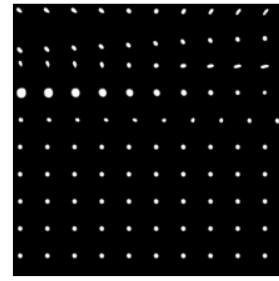
(c) Images traversal of dimensions



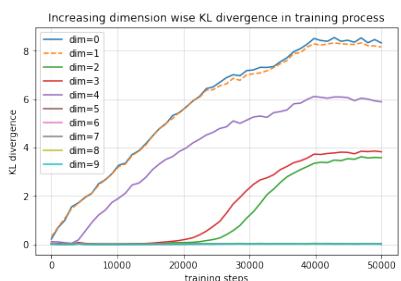
(d) Dimension wised KL in training



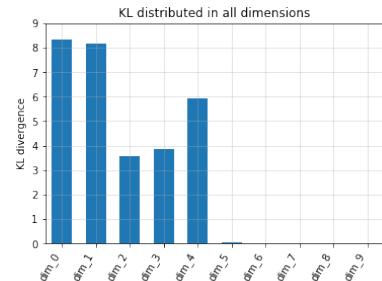
(e) Dimension wise KL in 50000 steps



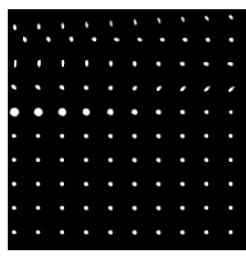
(f) Images traversal of dimensions



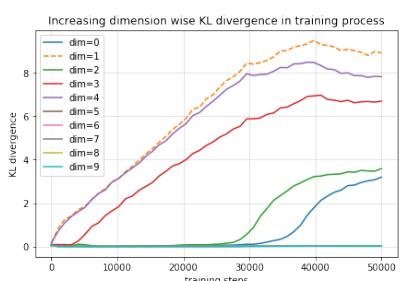
(a) Dimension wised KL in training



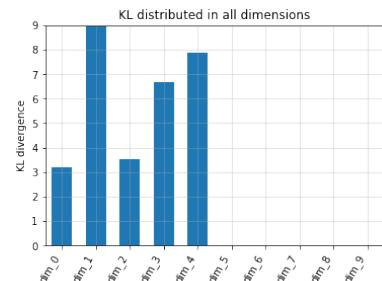
(b) Dimension wise KL in 50000 steps



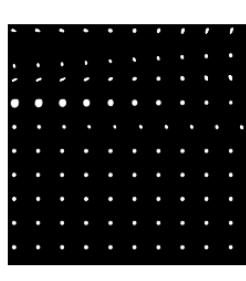
(c) Images traversal of dimensions



(d) Dimension wised KL in training



(e) Dimension wise KL in 50000 steps



(f) Images traversal of dimensions

Figure 20. Meaningful dimensions visualization for  $C_{max} = 30$ , end = 40000 (different random seeds). The KL divergences of target dimension (0-4 dimension) increase one by one during training (a). The KL divergences in different dimensions are different amounts after training (b). As the the image traversal results (c) shows, the meaningful dimensions are learned in 0-4 dims. So as to (d), (e) and (f).

## J. More Latent Space Visualizations

To illustrate that groupifying suppresses the latent space collapse in VAE-based methods is not an exception, we provide more visualizations of the latent space visualization of groupified VAEs. The original and groupified  $\beta$ -VAEs, AnnealVAEs, FactorVAEs and  $\beta$ -TCVAEs are shown in Figure 22, Figure 23, Figure 24 and Figure 25 respectively, which implies that the latent space of groupified VAEs are organized better than the original ones.

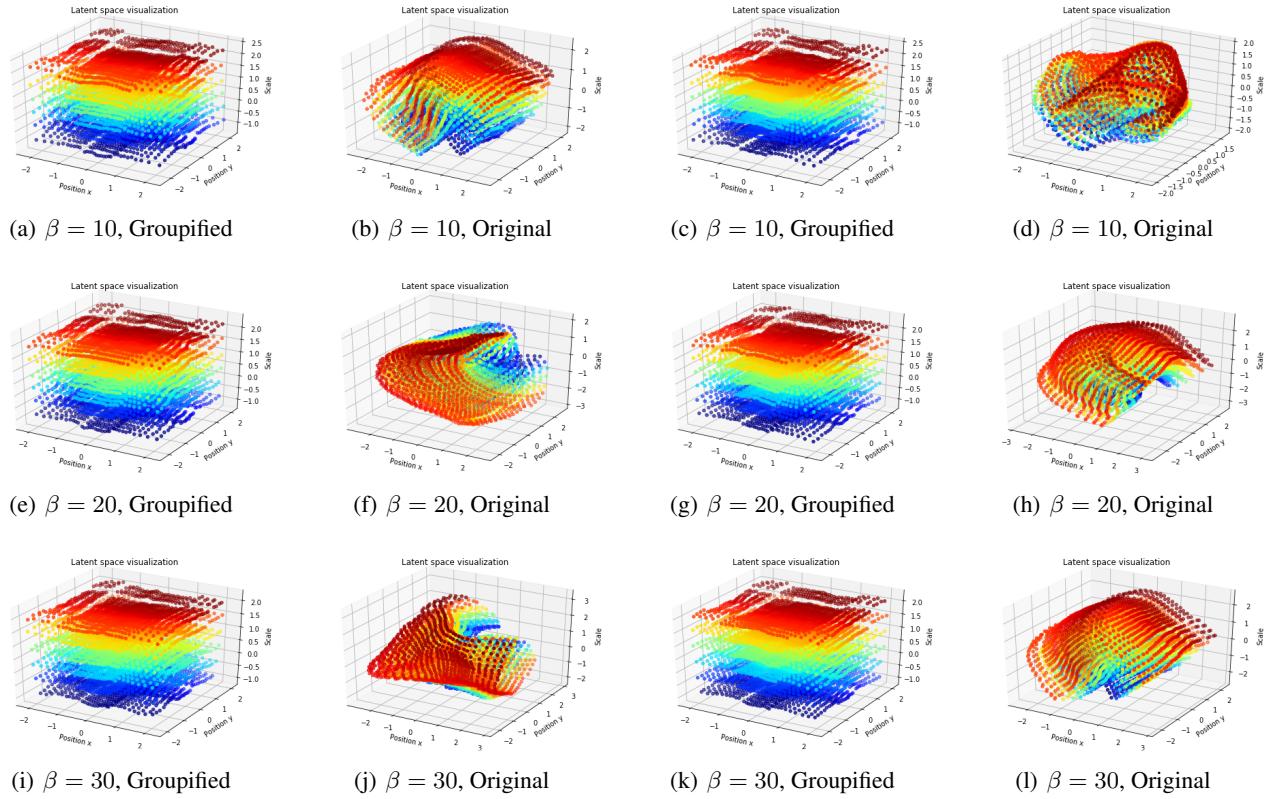
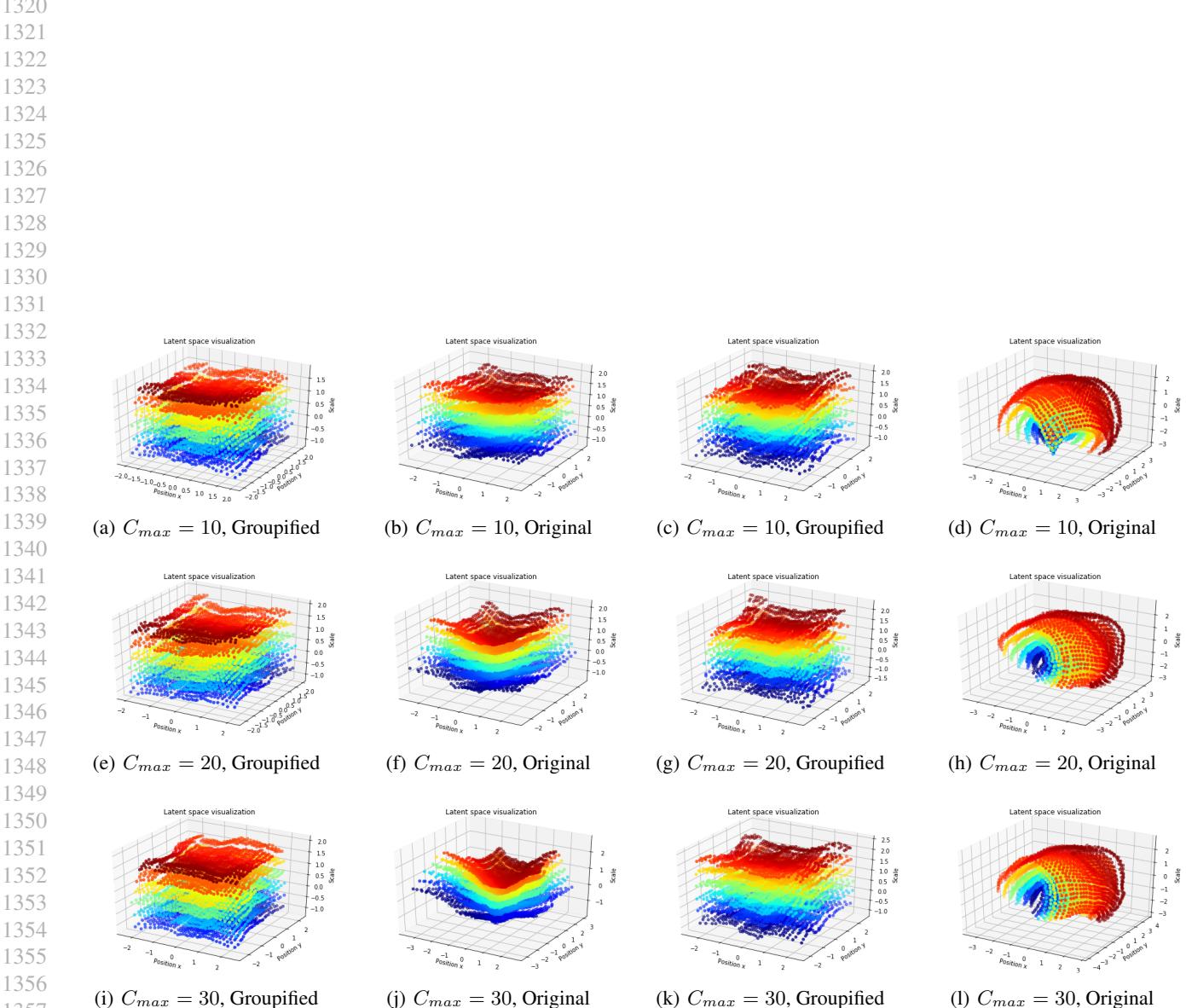
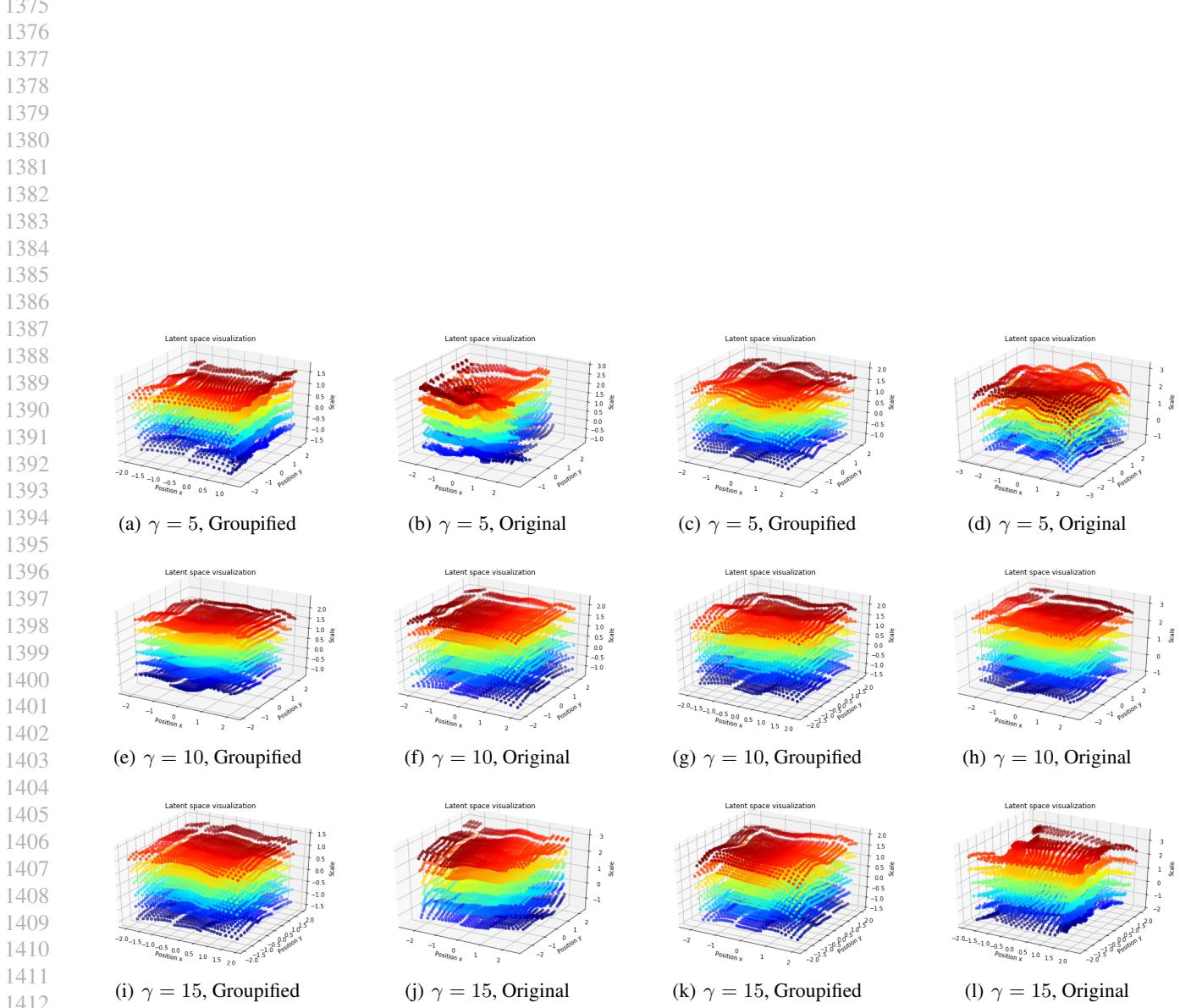


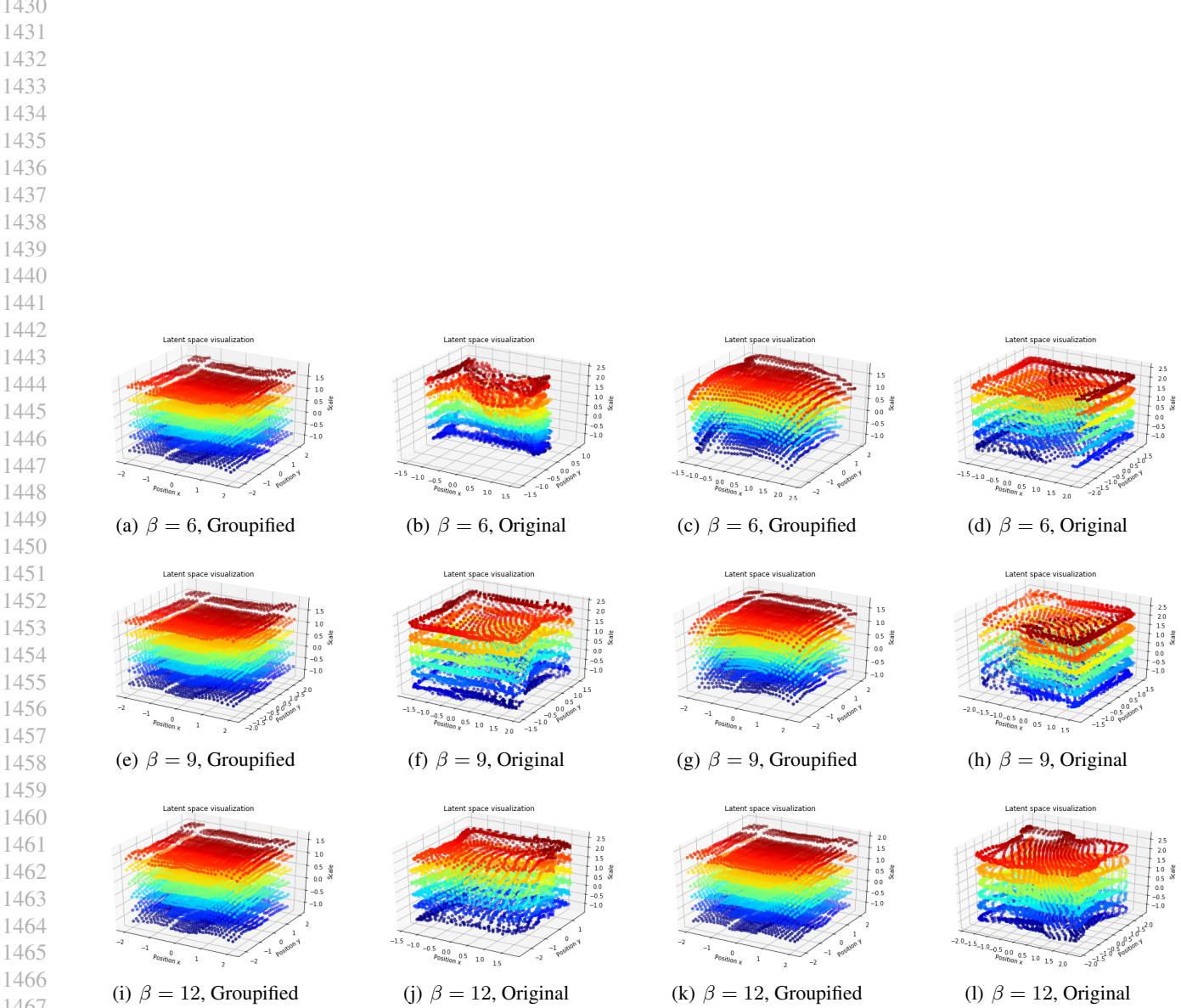
Figure 22. The latent space span by the groupified and original  $\beta$ -VAE. The models with same hyperparameter are trianed with different random seeds. The localization of each point is the disentangled representation of the corresponding image. And an ideal result is all the points form a cube and color variation is continuous. Higher hyper-parameter  $\beta$  results collapse of latent space. The collapse is suppressed by groupifying which lead to better disentanglement.



*Figure 23.* The latent space span by the groupified and original AnnealVAE. The models with same hyperparameter are trianed with different random seeds. The localization of each point is the disentangled representation of the corresponding image. And an ideal result is all the points form a cube and color variation is continuous. Higher hyper-parameter  $C_{max}$  results collapse of latent space. The collapse is suppressed by groupifying which lead to better disentanglement.



*Figure 24.* The latent space span by the groupified and original FactorVAE. The models with same hyperparameter are trianed with different random seeds. The localization of each point is the disentangled representation of the corresponding image. And an ideal result is all the points form a cube and color variation is continuous. The collapse of the latent space is suppressed by groupifying which lead to better disentanglement.



*Figure 25.* The latent space span by the groupified and original  $\beta$ -TCVAE. The models with same hyperparameter are trianed with different random seeds. The localization of each point is the disentangled representation of the corresponding image. And an ideal result is all the points form a cube and color variation is continuous. The collapse of the latent space is suppressed by groupifying which lead to better disentanglement.

**References**

1485  
1486 Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. beta-vae: Learning  
1487 basic visual concepts with a constrained variational framework. 2016.

1488  
1489 Hintze, J. L. and Nelson, R. D. Violin plots: a box plot-density trace synergism. *The American Statistician*, 52(2):181–184,  
1490 1998.

1491  
1492 Kim, H. and Mnih, A. Disentangling by factorising. 2018.

1493  
1494 Locatello, F., Bauer, S., Lucic, M., Raetsch, G., Gelly, S., Schölkopf, B., and Bachem, O. Challenging common assumptions  
1495 in the unsupervised learning of disentangled representations. In *ICML*, pp. 4114–4124, 2019.

1496  
1497 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A.  
1498 Automatic differentiation in pytorch. 2017.

1499  
1500 Reed, S. E., Zhang, Y., Zhang, Y., and Lee, H. Deep visual analogy-making. In *NeurIPS*, 2015.

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539