

# The Impact of Enforcing Representational Consistency of Identical Transformations for Disentangled Representation

Anonymous authors  
Paper under double-blind review

## Abstract

Recent symmetry-based approaches in Variational Autoencoders (VAEs) have advanced disentanglement learning and compositional generalization. However, existing methods can encode identical semantic transformations differently depending on the specific sample pairs, which reduce the *representational consistency of identical transformations*. In this paper, we analyze how three commonly used symmetry parameterization families in prior work, namely (1) matrix-exponential parameterizations over the general linear group  $GL(n)$ , (2) vector-additive actions in latent space, and (3) surjective mappings from latent vectors to the unit circle, can make it difficult to represent identical transformations consistently in dimension-wise disentangled latent spaces. To address this issue, we propose a framework that maps latent vectors to a bijective cyclic representation on the unit circle via the Cayley transform, together with a fixed-grid codebook regularization. We study this problem in a controlled setting and develop practical weakly supervised and supervised variants. Experiments on disentanglement benchmarks and compositional generalization tasks show that the proposed framework yields improved disentanglement performance and strong compositional generalization under supervised settings, with the stronger-supervision variants providing empirical reference points for the representational capacity of the framework. Overall, our results suggest that consistent representation of identical transformations is a useful design principle for improving disentanglement and generalization performance in the considered setting.

## 1 Introduction

Disentangled representation learning (Bengio et al., 2013; Wang et al., 2023) is a central problem in representation learning, as factorized latent representations can improve generalization under novel contexts and support more reliable reasoning over the underlying factors of variation (Montero et al., 2021). Since Locatello et al. (2019) establish that unsupervised disentanglement learning is impossible without suitable inductive bias (Bengio et al., 2013; Higgins et al., 2018), a major line of research focuses on introducing inductive biases that encourage latent variables to align with the factors of variation.

Among such approaches, symmetry-based modeling emerge as a promising direction. In particular, several recent works incorporate equivariant function modeling into Variational AutoEncoders (VAEs) by explicitly defining symmetry groups and their actions in latent space (Higgins et al., 2018). Concretely, prior methods instantiate specific symmetry parameterizations, including matrix Lie groups with matrix multiplication (Zhu et al., 2021; Keurti et al., 2023; Jung et al., 2024; Winter et al., 2022), cyclic-group formulations induced by surjective projections to the unit circle (Yang et al., 2021; Tonnaer et al., 2022; Cha & Thiyagalingam, 2023), and vector-additive constructions (Balabin et al., 2024). These studies show that symmetry-aware equivariant modeling improves disentanglement without supervision, thereby highlighting the usefulness of symmetry representations as an inductive bias.

However, an important issue remains insufficiently addressed: even when the underlying transformation is identical, existing approaches encode that transformation differently depending on which pair of samples induces it (Hwang et al., 2023). Such pair-dependent symmetry representations are undesirable for disentanglement

learning, because disentangled latent spaces are expected to capture the same factor change in a stable manner across different instances (Higgins et al., 2018). More broadly, this behavior contrasts with how humans tend to conceptualize transformations: the same transformation is often understood as a content-independent rule, or relational schema, that can be applied consistently across different contexts (Marcus et al., 1999; Gentner, 1983). Motivated by this perspective, we argue that an effective symmetry representation for disentanglement should encode the identical transformations consistently, independently of the particular sample pair from which it is inferred.

This consistency issue naturally raises a more fundamental question: what form of symmetry representation is appropriate for learning disentangled representations? Although prior works show the usefulness of symmetry-based inductive bias, they offer limited analysis of which symmetry parameterization is most suitable for consistently encoding identical transformations. Moreover, existing approaches often adopt a specific representation without thoroughly comparing how alternative parameterizations influence the consistency and suitability of the latent symmetry representation. We argue that this missing analysis is important because the choice of parameterization can directly influence whether identical transformations are represented consistently in latent space. Since our goal is to isolate and evaluate the role of symmetry parameterization itself, we study this question in a supervised setting, which provides a more controlled testbed than a fully unsupervised formulation and serves as a necessary step toward future unsupervised extensions.

In this work, we address these issues jointly: (1) the lack of mechanisms that enforce representational consistency of identical transformations, and (2) the limited analysis of which symmetry parameterizations are appropriate for achieving such consistency in disentanglement learning. To this end, we first analyze three commonly used symmetry parameterization families in prior work: matrix-exponential parameterizations over the general linear group (Zhu et al., 2021; Keurti et al., 2023; Jung et al., 2024; 2026; Winter et al., 2022), vector-additive actions in latent space (Balabin et al., 2024; Hwang et al., 2023), and surjective mappings from latent vectors to the unit circle (Yang et al., 2021; Tonnaer et al., 2022; Cha & Thiyagalingam, 2023). We show that these formulations have intrinsic limitations in representing identical transformations consistently in disentangled latent spaces. Building on this analysis, we propose a bijective representation on the unit circle based on the Cayley transform (Kreyszig et al., 2011), which provides a more suitable basis for consistent symmetry encoding. In practice, we further realize this principle through a fixed grid (codebook) (Hsu et al., 2023), which guides the model to represent identical transformations in latent space.

Our main contributions are as follows:

1. We identify a previously underexplored issue in symmetry-based disentanglement learning, the representational inconsistency of identical transformations across different sample pairs, and theoretically show that this issue is closely tied to the choice of symmetry parameterization by analyzing three widely used parameterization families from prior work.
2. We propose a bijective symmetry representation on the unit circle, together with a codebook-based latent vector space, to guide the model toward representational consistency of identical transformations.
3. We develop a practical learning framework that induces equivariance while promoting symmetry representation consistency, and empirically evaluate its effects on disentanglement learning and compositional generalization controlled benchmarks.

## 2 Preliminaries: Group Theory

**Binary operation:** Binary operation on a set  $S$  is a function that  $*$ :  $S \times S \rightarrow S$ , where  $\times$  is a cartesian product.

**Group:** A group is a set  $G$  together with binary operation  $*$ , that combines any two elements  $g_a$  and  $g_b$  in  $G$ , such that the following properties:

- closure:  $g_a, g_b \in G \Rightarrow g_a * g_b \in G$ .
- Associativity:  $\forall g_a, g_b, g_c \in G, s.t. (g_a * g_b) * g_c = g_a * (g_b * g_c)$ .

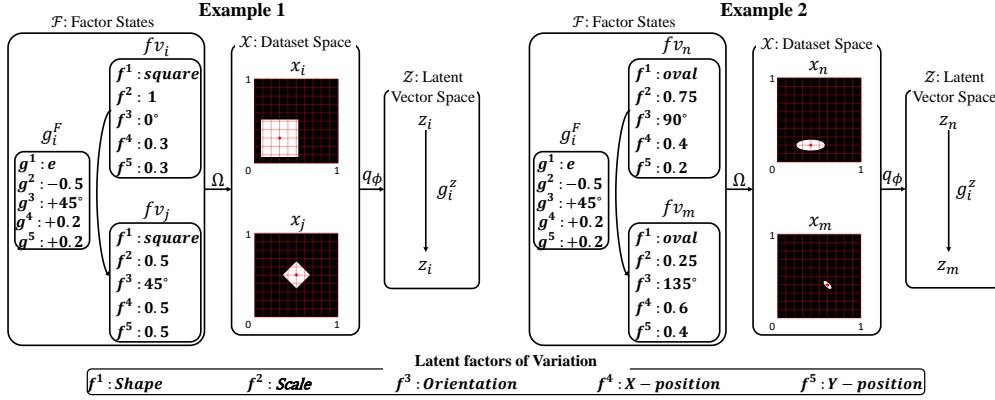


Figure 1: Representational consistency of identical transformations between example 1 and example 2 in mapping factor states, datasets, and latent representations. Equivariant function:  $q_\phi \circ \Omega$  satisfies  $g_i^z \cdot q_\phi \circ \Omega(fv_i) = q_\phi \circ \Omega(g_i^F \cdot fv_i)$ . Representational consistency of identical transformations: the same symmetry  $g_i^F$  is acted on latent factors  $fv_i, fv_n$ , and the symmetry acted on the latent vector space is represented as a single symmetry  $g_i^z = \Gamma(g_i^F)$  regardless of pairs.

- Identity element: There exists an element  $e \in G$ , s.t.  $\forall g \in G, e * g = g * e = g$ .
- Inverse element:  $\forall g \in G, \exists g^{-1} \in G: g * g^{-1} = g^{-1} * g = e$ .

**Group action:** Let  $(G, *)$  be a group and set  $X$ , binary operation  $\cdot : G \times X \rightarrow X$ , such that following properties:

- Identity:  $e \cdot x = x$ , where  $e \in G, x \in X$ .
- Compatibility:  $\forall g_a, g_b \in G, x \in X, (g_a * g_b) \cdot x = g_a \cdot (g_b \cdot x)$ .

**Equivariant map:** Let  $G$  be a group and  $X_1, X_2$  be two sets with corresponding group action of  $G$  in each sets:  $T_g^{X_1}, T_g^{X_2}$ , where  $g \in G$ . Then a function  $f : X_1 \rightarrow X_2$  is equivariant if  $f(T_g^{X_1} \cdot X_1) = T_g^{X_2} \cdot f(X_1)$ .

**Group Isomorphism:** Given two groups  $(G, *)$  and  $(H, \cdot)$ , a group isomorphism is a bijective function  $f : G \rightarrow H$  that satisfies as follows:  $f(u * v) = f(u) \cdot f(v)$ , where  $\forall u, v \in G ((G, *) \cong (H, \cdot))$ . The factor group  $\mathbb{R}/\mathbb{Z}$  and the circle group  $S^1$  are isomorphic, where group  $\mathbb{Z}$  of integers with addition, and  $S^1 = \{z \in \mathbb{C} : |z| = 1\}$ . ( $\mathbb{R}/\mathbb{Z} \cong S^1$ ). More details are in Appendix A.

### 3 Representational Consistency of Identical Transformations

#### 3.1 Equivariant Map from Factor States to Latent Vector Space

**Factor States.** The factor (world) states is referred to from Higgins et al. (2018) to define the equivariant function between the factor states and latent vector space for disentanglement learning with symmetries. We then follow the factor stats as :

- Factor States  $\mathcal{F}$  composed with latent factors of variation  $fv = fv^1 \times fv^2 \times \dots \times fv^k$ , and  $F = F^1 \times F^2 \times \dots \times F^k$ , where  $fv^i \in F^i, fv \in F$ , and set of factors  $F \subset \mathcal{F}$  as shown in Higgins et al. (2018) and Figure 1.

**Symmetries in Factor States.**  $G^F = G^{F^1} \times G^{F^2} \times \dots \times G^{F^k}$  is acted on the factor states  $\mathcal{F}$ , e.g.)  $fv_j = g_i^F \circ fv_i$  as shown in Figure 1, where  $g_i^F \in G^F$ , and  $fv_i, fv_j \in F$ . Function  $\Omega$  maps factor states to dataset space  $\Omega : \mathcal{F} \rightarrow \mathcal{X}$  and function  $q_\phi$  maps dataset to latent vector space  $q_\phi : \mathcal{X} \rightarrow \mathcal{Z}$  as shown in Figure 1.

**Equivariant Function from Factor States to Latent Vector Space.** Modeling equivariant functions is a common approach for disentanglement learning (Zhu et al., 2021; Yang et al., 2021; Keurti et al., 2023; Jung et al., 2024) and compositional generalization (Hwang et al., 2023) to inject inductive bias. To address two tasks, we model the composite function  $q_\phi \circ \Omega$  to be an equivariant function  $q_\phi \circ \Omega : \mathcal{F} \rightarrow \mathcal{Z}$ . As follows the definition of equivariant function  $f$  in Section 2,  $q_\phi \circ \Omega$  satisfies  $g^z \cdot q_\phi \circ \Omega(fv) = q_\phi \circ \Omega(g^F \cdot fv)$ , where  $G^F, G^z$  are symmetry group acted on space  $\mathcal{F}, \mathcal{Z}$  respectively, and  $g^F \in G^F, g^z \in G^z$ .

### 3.2 What is the Representational Consistency of Identical Transformations?

We formalize the *representational consistency of identical transformations* as follows (Figure 1):

**Definition 3.1. Representational consistency of identical transformations.** Let  $(fv_i, fv_j)$  and  $(fv_n, fv_m)$  be two pairs of elements in  $\mathcal{F}$  such that

$$fv_j = g^F \cdot fv_i, \quad fv_m = g^F \cdot fv_n, \quad (1)$$

where  $g^F$  denotes the same symmetry acting on  $\mathcal{F}$ .

Let  $q_\phi \circ \Omega : \mathcal{F} \rightarrow \mathcal{Z}$  be a mapping from  $\mathcal{F}$  to a latent space  $\mathcal{Z}$ . The corresponding latent representations are given by

$$z_i = q_\phi \circ \Omega(fv_i), \quad z_j = q_\phi \circ \Omega(g^F \cdot fv_i), \quad z_n = q_\phi \circ \Omega(fv_n), \quad z_m = q_\phi \circ \Omega(g^F \cdot fv_n). \quad (2)$$

Suppose there exists a latent-space transformation  $\Gamma(g^F)$  such that

$$q_\phi \circ \Omega(g^F \cdot fv_i) = \Gamma(g^F) \cdot q_\phi \circ \Omega(fv_i), \quad (3)$$

$$q_\phi \circ \Omega(g^F \cdot fv_n) = \Gamma(g^F) \cdot q_\phi \circ \Omega(fv_n). \quad (4)$$

We say that  $g^F$  is represented *consistently* if it is mapped to the same latent transformation  $\Gamma(g^F)$  regardless of which source element in  $\mathcal{F}$  induces it. That is, the latent representation of the transformation depends only on  $g^F$ , and not on the particular pairwise context from which it is inferred.

**Limited Guarantee of Consistency.** Previous works (Zhu et al., 2021; Yang et al., 2021; Jung et al., 2024; Hwang et al., 2023) formulate equivariant mappings between the dataset space and the latent vector space. Although Hwang et al. (2023) shows improved coherence for the same transformation, the problem of consistently representing identical transformations still remains unresolved. This suggests that equivariance alone does not necessarily guarantee the consistent representation of identical transformations defined above. Therefore, we explicitly treat consistency as an additional requirement for learning disentangled representation, and discuss the limitations of previous formulations in Section 4.

## 4 Missing Consistency of Transformations in Disentanglement Learning

**Cyclic Structure of Latent Factors of Variation.** To make Definition 3.1 operational, we must specify a concrete symmetry group acting on the factor states. In standard disentanglement benchmarks (dSprites, 3D Shapes, MPI3D), the labels of latent factors of variation are represented as discrete integers in a finite range, e.g.,  $[0, k - 1]$ , and the same interval difference between labels corresponds to the same semantic change in the dataset space (e.g., a fixed rotation angle, a fixed position shift). When such a factor admits a modulo structure, its label space is naturally modeled by the cyclic group  $\mathbb{Z}_k$ . We therefore assume that such latent factors of variation are isomorphic to cyclic groups, and formulate the symmetry group acting on factor states as  $G^F = G^{F^1} \times G^{F^2} \times \dots \times G^{F^n}$ ,  $G^{F^i} \cong \mathbb{Z}_{|F^i|}$ , where  $|F^i| \in \mathbb{Z}^+$  is the number of distinct values of the  $i$ -th factor. This assumption holds for all datasets used in this paper.

**Lack of Analysis on Suitable Symmetry Parameterizations for Disentangled Representations.** The goal of dimension-wise disentangled representations (Bengio et al., 2013; Wang et al., 2023) with symmetries ( $g_z \in G_z$ ) is to allow each symmetry to change a single latent dimension value. However, there is

limited analysis of which symmetry parameterization families are suitable for achieving a representational consistency of identical transformations in disentangled representations. In this section, we analyze the limitations of three commonly used symmetry parameterization families with proof sketches. In particular, we show that only the identity element yields a representational consistency of identical transformations in cases 1 and 2, while symmetry information is not preserved in case 3.

**Case 1: A Limitation of Matrix-Exponential Parameterizations over  $GL(n)$ .** Matrix-exponential parameterizations over the general linear group  $GL(n)$ , as used in prior works (Zhu et al., 2021; Kuzina et al., 2022; Miyato et al., 2022; Marchetti et al., 2023; Jung et al., 2026), are limited in achieving a representational consistency of identical transformations for disentangled representations. We first characterize the behavior of disentangled representations under matrix-exponential parameterizations, and then show the resulting limitation of  $GL(n)$ .

**Proposition 4.1.** *Let the symmetry group  $G_z$  ( $GL'(n)$ ) is defined as a subgroup of the General Linear group that implemented with matrix exponential, where  $GL'(n) = \{e^{\mathbf{M}} | \mathbf{M} \in \mathbb{R}^{n \times n}\}$ ,  $g^k$  is an element of  $GL'(n)$ , and  $g = \prod_k g^k$ . Then  $e^g \mathbf{z} = e \mathbf{I} g \mathbf{z} + \mathbf{v}'$ .*

**Theorem 4.2.** *(Limit of  $GL'(n)$ ) According to Proposition 4.1, only the identity matrix ( $g = \mathbf{I}$ ) represents the cyclic group of the dataset with representational consistency of identical transformations, where  $g \in GL'(n)$ .*

**Theorem 4.3.** *(Limit of  $GL(n)$ ) If  $H \subset GL(n)$ , then only the identity matrix of  $GL(n)$  represents the cyclic group of the dataset with representational consistency of identical transformations, where  $H = \{h | h = \mathbf{I} + \mathbf{M}^k\}$ ,  $\mathbf{m}^k$  is a column vector of  $\mathbf{M}^k$ ,  $\mathbf{m}^j = \vec{0}$  and  $j \in \{1, 2, \dots, n\} \setminus \{k\}$ .*

Therefore, the limitation of  $GL(n)$  is that only the  $\mathbf{I}$  represents the identical transformations with disentangled representation according to the Theorem 4.2, and Theorem 4.3. It implies that if  $g \neq \mathbf{I}$ , then  $g$  can not represent the identical transformations. More details are in Appendix B.2.

**Case 2: A Limitation of Vector-Additive Actions.** Another commonly used family is vector-additive actions in latent space, where the group action between two latent vectors is defined by vector addition, as in Balabin et al. (2024). We show that this family also has limitations in maintaining a representational consistency of identical transformations for cyclic groups.

**Corollary 4.4.** *If the group action is defined as  $\alpha(g, \mathbf{z}_i) = g + \mathbf{z}_i$ , then only zero vector represent the identical transformations for cyclic group with disentangled representation, where  $\mathbf{z} \in \mathbb{R}^n$ .*

According to the Corollary 4.4, it also shows a limitation in that only the identity element  $\vec{0}$  represents the identical transformations. More details of the proof are in Appendix B.3.

**Case 3: A Limitation of Surjective Mappings to the Unit Circle.** The third family consists of surjective mappings from latent vectors to the unit circle (Yang et al., 2021; Tonnaer et al., 2022; Cha & Thiyaalingam, 2023). We show that this family can cause undifferentiated symmetries under more general latent factors of variation and can lose part of the dataset’s symmetry information.

**Corollary 4.5.** *By the equivariant and surjective function  $b: \mathcal{Z} \rightarrow \mathcal{Y}$ , the capacity of  $\mathcal{Z}$  and  $\mathcal{Y}$  is  $|\mathcal{Z}| \geq |\mathcal{Y}|$  then  $\Gamma'$  is an endomorphism because  $|G_{\mathcal{Z}}| \geq |G_{\mathcal{Y}}|$ . On the other hand, isomorphism identically maps the two spaces ( $|G_{\mathcal{Z}}| = |G_{\mathcal{Y}}|$ ).*

Therefore, if  $b$  is surjective and not injective, then there exists at least one case where  $\Gamma'(g_i) = \Gamma'(g_j)$ . It implies that loss of symmetry structure occurs with a surjective function. More details of the proof are in Appendix B.4.

## 5 Method

### 5.1 Motivation: From Limitations to Bijective Cyclic Representation

The above analysis reveals two distinct types of limitation. Cases 1 and 2 exhibit a *Type I* failure: the consistent representation condition in Definition 3.1 is satisfied only by the trivial identity element, making

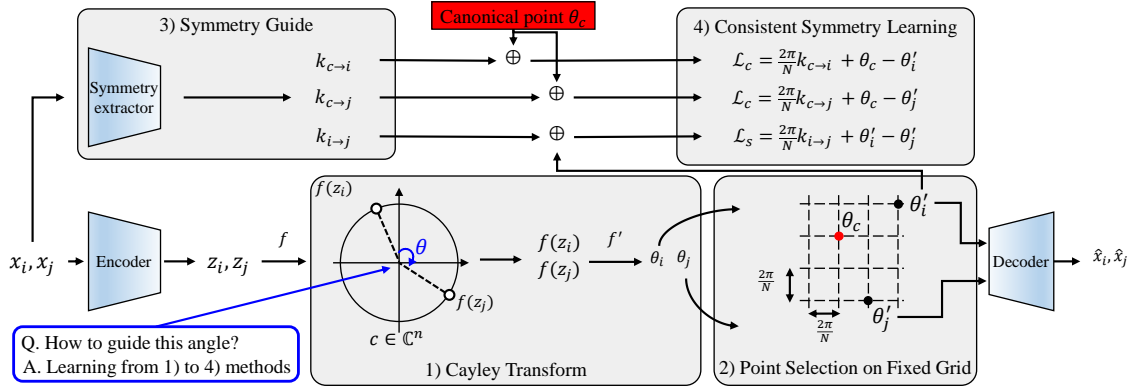


Figure 2: The overall architecture of our proposed method comprises four distinct components: 1) Cayley transform of latent vectors to angle space, 2) point selection of fixed grid for consistent symmetry, 3) defining the step size between two inputs through three methods, and 4) a loss function that satisfies the group action  $\alpha(g, \theta) = g + \theta$ .

non-trivial consistent encoding impossible. Case 3 exhibits a *Type II* failure: the surjective (but non-injective) mapping collapses distinct symmetry elements onto the same latent transformation, causing irreversible loss of group structure. Together, these results motivate a representation that is (i) bijective to preserve the full symmetry structure, and (ii) discretely structured so that identical transformations map to the same grid point regardless of the sample pair. We address both requirements in the following subsections, and formally show in Appendix C.1 that the resulting framework satisfies Definition 3.1 under the cyclic group structure.

## 5.2 Cayley Transform: Representational Consistency of Identical Transformations on Cyclic Representations

We first define the latent vector space as a  $G$ -set of a cyclic group to ensure the representational consistency of identical transformations because 1) we assume that  $G^{F^i}$  is a cyclic group as referred in Section 4 and 2) a single cyclic group element can represent all elements, as demonstrated in Higgins et al. (2018). To address the issues discussed in Section 4, we implement the Cayley transform (Kreyszig et al., 2011) that maps real numbers to complex numbers bijectively.

**Cyclic Group for Consistency.** As we assume that factor states serve as latent factors of variation for input samples in Section 4, the symmetry group  $G^F = \mathbb{Z}/|F^1|\mathbb{Z} \times \mathbb{Z}/|F^2|\mathbb{Z} \times \dots \times \mathbb{Z}/|F^n|\mathbb{Z}$ , where  $|F^k| \in \mathbb{Z}^+$  represents the number of factors in the datasets. Additionally, the cyclic group effectively represents the symmetries over the same group action, as the cyclic group  $G = \{e, g^1, g^2, \dots, g^{n-1}\}$  consists entirely of integer powers of the group element  $g$ . Therefore, if the model learns a single symmetry element  $g$ , it represents the entire symmetry group. Motivated by Yang et al. (2021), we implement the cyclic group as the  $n^{\text{th}}$  root of unity.

**Group Action and  $G^c$ -set.** As we define the cyclic group as  $n^{\text{th}}$  root of unity, the cyclic group  $G^c = G_1^c \times G_2^c \times \dots \times G_k^c$  and  $G_i^c = \{g_i^c | g_i^c = \frac{2\pi}{N}k, k \in \{0, 1, 2, \dots, N-1\}\}$ , where  $N \in \mathbb{Z}^+$ . We define the group action  $\alpha : \Theta^n \times G^c \rightarrow \Theta^n$  as  $\alpha(g^c, \theta) = g^c + \theta$ , where the latent vector  $\theta \in \Theta^n$ , and  $-\pi < \theta_i^k \leq \pi$ .

**Cayley Transform for Complex Number Space.** VAE frameworks establish the latent vector space in the real number space with the Gaussian normal distribution, so the latent vector space is not a  $G^c$ -set as we assume ( $z \in \Theta^n$ ). For the defined  $G^c$ -set, we utilize the Cayley transform and invertible function  $f' \circ f : [-\infty, \infty] \rightarrow \{\theta | -\pi < \theta \leq \pi\}$ . To map real numbers to the complex number space, we utilize the Cayley transform function (bijective)  $f : [-\infty, \infty] \rightarrow \{c_i^k \in \mathbb{C} : |c_i^k| = 1\}$ , defined as follows:

$$f(z_i^k) = c_i^k = \frac{z_i^k - i}{z_i^k + i} = i \frac{-2z_i^k}{(z_i^k)^2 + 1} + \frac{(z_i^k)^2 - 1}{(z_i^k)^2 + 1}, \quad (5)$$

Table 1: Comparison of symmetry representations with CTFG-GT models. We evaluate ground truth supervised models under different symmetry representations to investigate which representation is better suited to consistently encoding identical transformations. Specifically, we compare three symmetry modeling approaches—General Linear group ( $GL(n)$ ), vector addition (Add.), and surjective mappings (Sur.)—against our proposed method.

Symmetry	3D Shapes				dSprites			
	beta-VAE	FVM	MIG	DCI	beta-VAE	FVM	MIG	DCI
CTFG-GT ( $GL(n)$ )	73.50( $\pm 17.92$ )	66.16( $\pm 15.95$ )	19.66( $\pm 24.29$ )	37.93( $\pm 23.92$ )	65.11( $\pm 3.89$ )	47.69( $\pm 8.04$ )	2.12( $\pm 1.73$ )	5.90( $\pm 2.11$ )
CTFG-GT (Add.)	78.60( $\pm 11.43$ )	60.50( $\pm 15.29$ )	25.13( $\pm 17.12$ )	45.84( $\pm 8.36$ )	68.89( $\pm 11.45$ )	68.31( $\pm 9.48$ )	7.22( $\pm 4.31$ )	11.83( $\pm 3.85$ )
CTFG-GT (Sur.)	73.27( $\pm 9.97$ )	63.85( $\pm 6.81$ )	6.20( $\pm 4.03$ )	29.79( $\pm 9.86$ )	80.60( $\pm 10.83$ )	60.40( $\pm 6.47$ )	13.68( $\pm 3.86$ )	23.01( $\pm 2.28$ )
CTFG-GT	<b>100.00</b> ( $\pm 0.00$ )	<b>100.00</b> ( $\pm 0.00$ )	<b>96.57</b> ( $\pm 0.80$ )	<b>99.94</b> ( $\pm 0.18$ )	<b>95.80</b> ( $\pm 4.57$ )	<b>99.26</b> ( $\pm 1.12$ )	<b>51.81</b> ( $\pm 2.97$ )	<b>63.26</b> ( $\pm 2.73$ )

where the  $z_i^k$  is a  $k^{th}$  dimension value of  $z_i \in \mathbb{R}^n$ . We define a bijective function that maps complex numbers to the angle space for simplicity  $f' : \{c_i^k \in \mathbb{C} : |c_i^k| = 1\} \rightarrow \{\theta_i^k \mid -\pi < \theta_i^k \leq \pi\}$  as follows:

$$\theta_i^k = f'(c_i^k) = \begin{cases} \cos^{-1}(\Re(c_i^k)) - \pi, & \text{if } \Im(c_i^k) >= 0 \\ \pi - \cos^{-1}(\Re(c_i^k)) & \text{otherwise} \end{cases}, \quad (6)$$

where  $\Re(c_i^k)$  and  $\Im(c_i^k)$  are real and imaginary parts of  $c_i^k$ , respectively.

### 5.3 Point Selection on Fixed Grid for Representational Consistency of Identical Transformations

To address the cyclic group and consistency of identical transformations in the latent vector space, we set the space as a fixed grid (Mentzer et al., 2024) instead of a continuous or learnable grid (Hsu et al., 2023). Because, as shown in Figure 2, the interval between two nearest codes is always  $\frac{2\pi}{N_i}$ , and it implies that equation  $g_i^c = \frac{2\pi}{N_i}$  is satisfied for all cases. Then we utilized the finite scalar quantization (Mentzer et al., 2024) for fixed codebook  $\mathbf{V} \in \mathbb{R}^N$  as follows:

$$\mathbf{V} = [-\pi + \frac{2\pi}{N}, \dots, -\pi + \frac{2\pi}{N}k, \dots, -\pi + \frac{2\pi}{N}(N-1), \pi]. \quad (7)$$

Then we select the nearest neighbor of the latent vector as Hsu et al. (2023):  $\theta'^k = \arg \min_{\mathbf{V}^i \in \mathbf{V}} |\mathbf{V}^i - \theta^k|$ , where  $\mathbf{V}^i$  is the  $i^{th}$  dimension value of the codebook  $\mathbf{V}$ . We define the grid loss  $\mathcal{L}_{grid} = \|\theta' - \theta\|_2^2$  to consistently select the grid, where  $\|\cdot\|_2$  is a L2 norm.

### 5.4 Symmetry Guide: How to Select Step ( $k$ )?

As we define the cyclic group  $G_i^c = \{g_i^c \mid g_i^c = \frac{2\pi}{N_i}k, k \in \{0, 1, 2, \dots, N_i - 1\}\}$ , we implement the step  $k$  to guide how much step moves to be  $\theta_i = g^c + \theta_j$  (defined group action). We propose three guiding approaches: 1) ground truth, 2) supervised, 3) and weakly-supervised methods.

**Ground Truth Based Method.** As shown in Figure 2, we set the symmetry group elements  $g^c$  from the ground truth of samples as follows:

$$\mathbf{k}_{i \rightarrow j} = \begin{cases} l_j - l_i & \text{if } l_j - l_i \geq 0 \\ N + l_j - l_i & \text{otherwise} \end{cases}, \quad (8)$$

where  $\mathbf{k}_{i \rightarrow j}$  is a step size,  $g^c = \frac{2\pi}{N} \mathbf{k}_{i \rightarrow j}$ , and  $l_i$  and  $l_j$  are labels of samples  $x_i$  and  $x_j$ , respectively.

**Supervised Method.** As shown in Figure 2, we train the symmetry extractor to predict the labels of samples ( $\hat{l}$ ) using cross-entropy loss, defined as  $\mathcal{L}_{pred} = C.E.(\hat{l}, l)$ . We then define the symmetry group elements by  $\hat{l}$  instead of the ground truth labels  $l$ .

**Weakly-Supervised Method.** We utilize a  $p$  ratio of the labels for prediction, while the remaining labels are predicted using the pseudo-label loss as follows:  $\mathcal{L}_{pl} = \sum_i^{|F_i|} D_{KL}(p(l^i|x) \| p(l^i))$ , where the  $l^i$  is a  $i^{th}$  factor of the label and a discrete uniform distribution  $p(l^i) \sim \mathcal{U}\{1, |F_i|\}$ . We define the  $p(l^i|x)$  as the distribution of the classifier.

Table 2: Disentanglement performance on dSprites, 3D Shapes, and the MPI3D dataset. (Unsup: unsupervised method, Weak-sup: weakly supervised, Sup: supervised, GT: Ground-Truth, CTFG: our method). Bold text indicates a higher value than the other baseline models.

type	Method	dSprites				3D Shapes				MPI3D			
		beta-VAE	FVM	MIG	DCI	beta-VAE	FVM	MIG	DCI	beta-VAE	FVM	MIG	DCI
Unsup	$\beta$ -VAE	78.40( $\pm 9.03$ )	64.84( $\pm 11.40$ )	14.52( $\pm 9.33$ )	22.37( $\pm 11.80$ )	90.33( $\pm 7.42$ )	72.63( $\pm 19.55$ )	40.49( $\pm 23.31$ )	54.32( $\pm 16.45$ )	57.60( $\pm 7.93$ )	40.86( $\pm 3.92$ )	4.91( $\pm 1.43$ )	22.29( $\pm 1.42$ )
	$\beta$ -TCVAE	81.80( $\pm 11.91$ )	70.16( $\pm 12.41$ )	19.03( $\pm 9.40$ )	30.89( $\pm 8.96$ )	88.20( $\pm 7.91$ )	74.45( $\pm 14.68$ )	43.17( $\pm 28.28$ )	59.71( $\pm 14.79$ )	55.40( $\pm 9.52$ )	40.80( $\pm 2.60$ )	5.23( $\pm 1.96$ )	21.47( $\pm 2.35$ )
	Factor-VAE	87.20( $\pm 7.50$ )	76.80( $\pm 7.50$ )	24.98( $\pm 12.03$ )	33.38( $\pm 12.27$ )	95.60( $\pm 6.99$ )	80.53( $\pm 11.41$ )	54.54( $\pm 25.10$ )	66.55( $\pm 9.52$ )	54.00( $\pm 7.18$ )	39.64( $\pm 3.81$ )	4.34( $\pm 0.69$ )	21.24( $\pm 2.04$ )
	CLG-VAE	88.40( $\pm 5.80$ )	82.21( $\pm 5.73$ )	20.89( $\pm 7.40$ )	29.96( $\pm 7.05$ )	86.20( $\pm 5.61$ )	77.36( $\pm 7.99$ )	50.39( $\pm 12.37$ )	59.25( $\pm 11.21$ )	46.40( $\pm 7.35$ )	37.31( $\pm 2.27$ )	20.77( $\pm 5.70$ )	24.26( $\pm 2.73$ )
Weak-sup	Ada-GVAE	83.60( $\pm 2.61$ )	83.67( $\pm 2.97$ )	21.34( $\pm 5.35$ )	47.26( $\pm 1.89$ )	72.75( $\pm 6.50$ )	59.81( $\pm 6.14$ )	24.77( $\pm 7.48$ )	64.57( $\pm 4.04$ )	64.89( $\pm 7.22$ )	46.10( $\pm 3.19$ )	22.48( $\pm 8.14$ )	41.30( $\pm 7.00$ )
	CTFG-wSP	<b>90.50(<math>\pm 7.55</math>)</b>	<b>84.50(<math>\pm 1.41</math>)</b>	<b>31.95(<math>\pm 2.40</math>)</b>	39.36( $\pm 1.49$ )	<b>95.00(<math>\pm 7.07</math>)</b>	<b>88.81(<math>\pm 13.17</math>)</b>	<b>57.94(<math>\pm 16.52</math>)</b>	<b>72.14(<math>\pm 3.23</math>)</b>	<b>67.50(<math>\pm 12.68</math>)</b>	<b>81.59(<math>\pm 3.80</math>)</b>	<b>61.07(<math>\pm 4.81</math>)</b>	<b>78.15(<math>\pm 0.57</math>)</b>
Sup	CTFG-SP	<b>91.40(<math>\pm 4.99</math>)</b>	<b>93.74(<math>\pm 1.82</math>)</b>	<b>51.02(<math>\pm 2.42</math>)</b>	<b>64.69(<math>\pm 1.55</math>)</b>	<b>100.00(<math>\pm 0.00</math>)</b>	<b>100.00(<math>\pm 0.00</math>)</b>	<b>96.95(<math>\pm 0.18</math>)</b>	<b>99.99(<math>\pm 0.01</math>)</b>	<b>86.40(<math>\pm 8.63</math>)</b>	<b>99.96(<math>\pm 0.08</math>)</b>	<b>62.78(<math>\pm 6.95</math>)</b>	<b>88.06(<math>\pm 1.66</math>)</b>
GT	CTFG-GT	<b>95.80(<math>\pm 4.57</math>)</b>	<b>99.26(<math>\pm 1.12</math>)</b>	<b>51.81(<math>\pm 2.97</math>)</b>	<b>63.26(<math>\pm 2.73</math>)</b>	<b>100.00(<math>\pm 0.00</math>)</b>	<b>100.00(<math>\pm 0.00</math>)</b>	<b>96.57(<math>\pm 0.80</math>)</b>	<b>99.94(<math>\pm 0.18</math>)</b>	<b>76.80(<math>\pm 9.66</math>)</b>	<b>99.03(<math>\pm 0.98</math>)</b>	<b>65.17(<math>\pm 8.11</math>)</b>	<b>81.12(<math>\pm 1.03</math>)</b>

## 5.5 Objective Function for Representational Consistency of Identical Transformations

The symmetry loss below enforces relative consistency between a given pair  $(x_i, x_j)$ , but the induced latent transformation still depends on the specific pair, which is precisely the pair-dependency issue in Definition 3.1. To eliminate this dependency, we additionally introduce a learnable canonical reference point  $\theta_c$  that anchors each sample to an absolute position in latent space independently of which pair induced the transformation.

**Symmetry Loss (Relative Position).** As we define the group action  $\alpha(g^c, \theta_i)$  and step size  $\mathbf{k}$ , we implement the symmetry loss  $\mathcal{L}_s$  to satisfy  $\theta_j = \frac{2\pi}{N}\mathbf{k}_{i \rightarrow j} + \theta_i$ :

$$\mathcal{L}_s = \|f' \circ f \circ q_\phi(x_j) - \left(\frac{2\pi}{N}\mathbf{k}_{i \rightarrow j} + f' \circ f \circ q_\phi(x_i)\right)\|_2^2. \quad (9)$$

We set the code loss  $\mathcal{L}_{code} = \mathcal{L}_{grid} + \mathcal{L}_s$ .

**Canonical Loss (Absolute Position)** The defined symmetry  $g^c$  represents the movement between two latent vectors  $(\theta_i$  and  $\theta_j)$ . It implies that learning symmetries depends on pairs of observations. To eliminate this dependency on specific observations, we propose learning absolute transformations through a learnable canonical point  $\theta_c$  to satisfy  $\theta_i = \frac{2\pi}{N}\mathbf{k}_{c \rightarrow i} + \theta_c$ :

$$\mathcal{L}_c = \|f' \circ f \circ q_\phi(x_i) - \left(\frac{2\pi}{N}\mathbf{k}_{c \rightarrow i} + \theta_c\right)\|_1, \quad (10)$$

where  $\theta_c$  is a learnable canonical point  $\theta_c \in \Theta^n$ ,  $\mathbf{k}_{c \rightarrow i} = l_i$ , and  $\|\cdot\|_1$  is a L1 norm.

**Objective Loss.** Our objective losses are defined as 1)  $\mathcal{L}_{GT} = \mathcal{L}_{recont} + \alpha\mathcal{L}_{code} + \gamma\mathcal{L}_c$  for Cayley transform and fixed grid Ground Truth model (CTFG-GT), 2)  $\mathcal{L}_{Sup} = \mathcal{L}_{GT} + \beta\mathcal{L}_{pred}$  for supervised method (CTFG-SP), and 3)  $\mathcal{L}_{weak-Sup} = \mathcal{L}_{Sup} + \lambda\mathcal{L}_{pl}$  for weakly-supervised method (CTFG-wSP).

## 6 Experiments

**Common Datasets.** We utilize the dSprites (Matthey et al., 2017), 3D Shapes (Burgess & Kim, 2018), and MPI3D (Gondal et al., 2019) datasets for compositional generalization and disentanglement learning tasks. More details are in Appendix D.2.

### 6.1 Disentanglement Learning

**Settings.** We set the common hyper-parameters of the proposed method  $\alpha \in \{100, 1000\}$ ,  $\gamma = 1$  for supervised and ground truth model,  $\beta \in \{1.0, 2.0\}$  for supervised method, and  $\lambda = 1.0, p = 0.5$  for weakly-supervised method. We run 10 seed variance over each model with seed  $\in \{1, 2, \dots, 10\}$ . More details are in Appendix D.4.

**Case Studies: Do Any Ground Truth Based Methods Encourage Representational Consistency of Identical Transformations?** As shown in Section 4, we briefly show the difficulty of prior symmetry representations ( $GL(n)$ , vector addition, and subjective function) for the disentangled representation following

ground truth based method in Section 5.4. We demonstrate that previous methods are limited in preserving the dataset’s symmetry structure. Consequently, we have adopted these three types of symmetry instead of our method. As indicated in Table 1, the CTFG-GT method outperforms other methods across all metrics. This suggests that enforcing the representational consistency of identical transformations in the angle space is a more suitable method for disentanglement learning.

**Quantitative Results.** Although comparisons with unsupervised methods should be interpreted cautiously due to the different supervision regimes, CTFG-SP and CTFG-GT achieve consistently strong performance on datasets, providing empirical reference points for the representational capacity of the framework, as shown in Table 2. In particular, they achieve near-saturated disentanglement scores on 3D Shapes and consistently strong performance on dSprites and MPI3D, indicating that the proposed framework remains effective under stronger supervision in these benchmark settings.

More importantly, under a direct weakly supervised comparison, CTFG-wSP generally outperforms Ada-GVAE, achieving higher scores on 11 out of 12 dataset-metric pairs. These results support the effectiveness of our method in learning disentangled representations under comparable weak supervision.

**Disentanglement vs. Reconstruction.** Most disentangled representation learning models face a trade-off between reconstruction error and disentanglement metrics (Kingma & Welling, 2013; Chen et al., 2018; Kim & Mnih, 2018; Zhu et al., 2021; Keurti et al., 2023; Higgins et al., 2017; Locatello et al., 2020). However, our results suggest that this trade-off can be mitigated in the evaluated setting, as illustrated in Figure 3. Although our model’s reconstruction error is two times lower than the baselines, it achieves higher disentanglement performance than the others with the MPI3D datasets. Further details are provided in Appendix F.1.

### Direct Measurement of Symmetry Consistency.

To directly validate Definition 3.1, we measure the cosine similarity between latent differences induced by identical ground-truth transformations. Specifically, for 50,000 sample pairs satisfying  $l_1 - l_2 = l_3 - l_4$  under ground-truth labels, we compute  $\cos(\theta_1 - \theta_2, \theta_3 - \theta_4)$  over 10 random seeds. As shown in Table 3, all CTFG variants achieve near-perfect consistency across all datasets: CTFG-GT attains 1.0000 on every benchmark, while CTFG-SP and CTFG-wSP reach above 0.99 and 0.96, respectively. These results confirm that the proposed framework produces stable latent representations of identical transformations across different sample pairs, directly supporting the claim of Definition 3.1.

Table 3: Symmetry consistency scores across datasets.

	CTFG-sWP	CTFG-SP	CTFG-GT
dSprites	0.9693( $\pm 0.0002$ )	0.9889( $\pm 0.0001$ )	1.0000( $\pm 0.0000$ )
3D Shapes	0.9835( $\pm 0.0002$ )	0.9982( $\pm 0.0001$ )	1.0000( $\pm 0.0000$ )
MPI3D	0.9853( $\pm 0.0001$ )	0.9984( $\pm 0.0001$ )	1.0000( $\pm 0.0000$ )

**Qualitative Analysis.** As shown in Figure 4a–4f, the baseline results show that multiple factors are changed when a single dimension value is changed on both the 3D Shapes and MPI3D datasets. Also, objects disappear at certain intervals in the baseline results. On the other hand, CTFG-SP and CTFG-GT show better results than the baselines. Compared to the baseline model, the cases of overlapping factors in a single dimension are reduced by the proposed models on both the 3D Shapes and MPI3D datasets.

## 6.2 Compositional Generalization

As shown in Section 6.1, encouraging a representational consistency of identical transformations leads to improved disentanglement performance in controlled settings. Because an important goal of disentangled

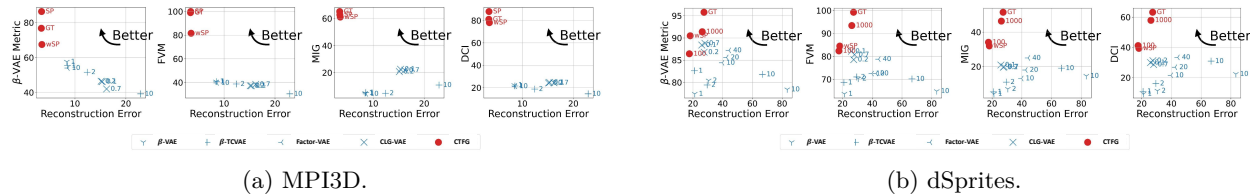


Figure 3: Reconstruction error vs. evaluation metrics of the MPI3D and dSprites dataset ( $\beta$ -VAE metric, FVM, MIG, and DCI). The top left side indicates the best results on both objectives

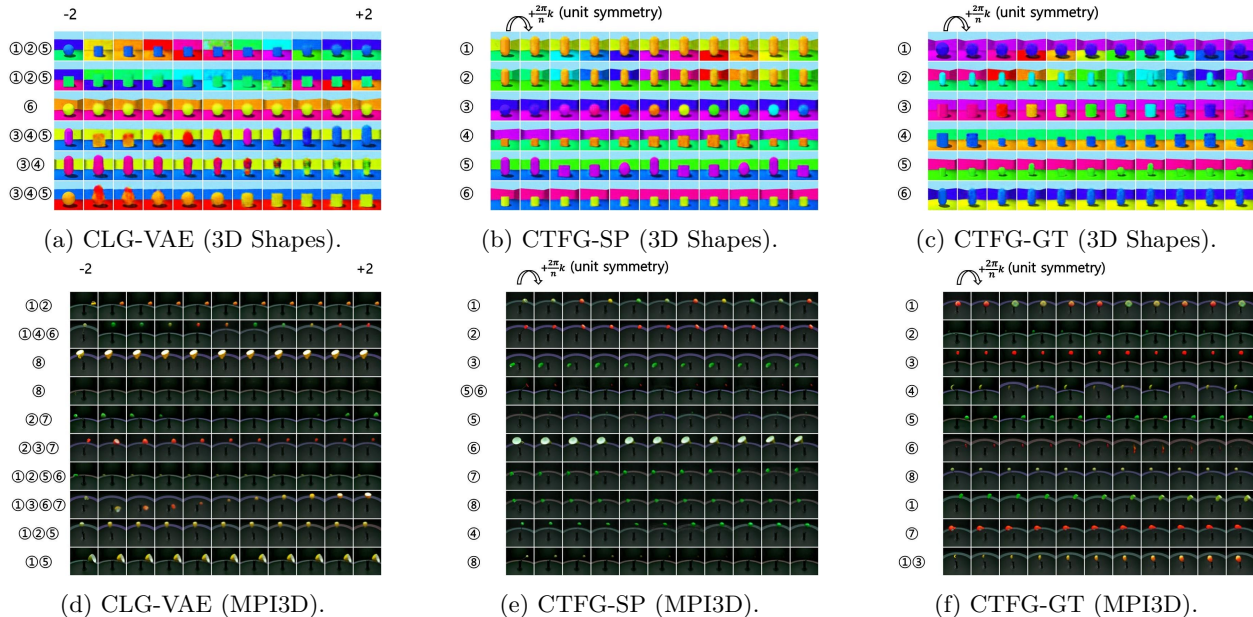


Figure 4: Alignment of a factor and a dimension: How many factors are changed following the dimension (column-wise direction)? The 1<sup>st</sup> column images are randomly selected from the dataset. Each row indicates each dimension of each model. The Commutative Lie Group VAE trace each dimension value from -2 to +2. The proposed methods apply a group action  $+\frac{2\pi}{n}$  to the selected images a total of 10 times. The numbers in Figure 4a-4c refer to factors of the 3D Shapes dataset: ①, ②, ③, ④, ⑤, and ⑥ refer to floor color, wall color, object color, scale, shape, and orientation, respectively. The numbers in Figure 4d-4f refer to factors of the MPI3D dataset: ①, ②, ③, ④, ⑤, ⑥, and ⑦ refer to object color, object shape, object size, camera height, background color, horizontal axis, and vertical axis, respectively.

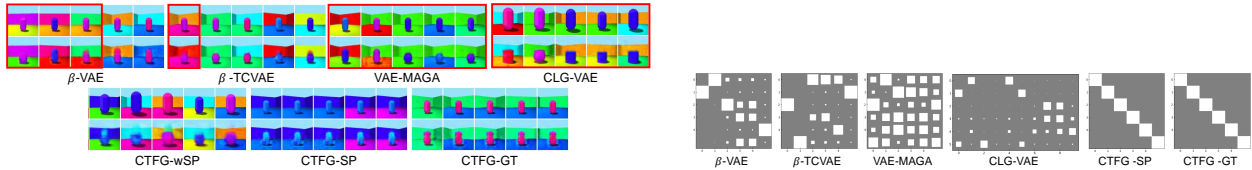
representations is to support generalization, we further evaluate our approach in compositional generalization settings.

**Settings.** We excess Recombination-to-Element (R2E) and the Recombination-to-Range (R2R) tasks. We separate the training and test datasets following previous studies (Montero et al., 2021; Hwang et al., 2023), with additional details and hyper-parameter tuning provided in Appendix D.3. We run each model with three seeds  $\in \{1, 2, 3\}$ . We assess the reconstruction error, a general compositional generalization metric.

**Quantitative Results.** As shown in Table 4, the proposed method CTFG-GT is significantly improved with all datasets. It implies that the representational consistency of identical transformations also impacts to compositional generalization. Also, the supervised method CTFG-SP demonstrates advancements in all

Table 4: Compositional generalization performance of dSprites, 3D Shapes, and MPI3D datasets. We select the best results from the hyper-parameter tunings. The evaluation metric is the reconstruction error (BCE loss for dSprites, and MSE loss for 3D Shapes and MPI3D).

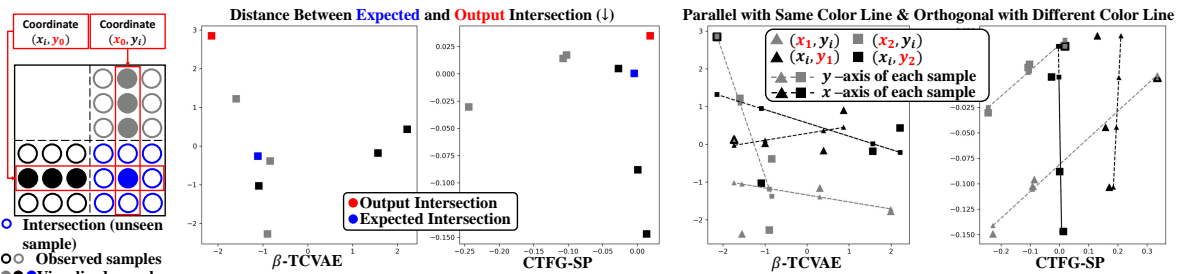
Method	dSprites		3D Shapes		MPI3D		
	R2E	R2R	R2E	R2R	R2E	R2R	
Base	$\beta$ -VAE	10.85( $\pm 0.67$ )	179.52( $\pm 12.15$ )	16.59( $\pm 1.72$ )	268.59( $\pm 76.59$ )	6.63( $\pm 0.65$ )	8.50( $\pm 0.55$ )
	$\beta$ -TCVAE	10.73( $\pm 0.03$ )	153.75( $\pm 7.65$ )	14.74( $\pm 0.14$ )	221.72( $\pm 41.57$ )	5.60( $\pm 0.21$ )	8.73( $\pm 0.36$ )
Equiv.	CLG-VAE	10.30( $\pm 0.20$ )	246.37( $\pm 53.42$ )	21.86( $\pm 0.98$ )	276.81( $\pm 17.01$ )	11.68( $\pm 0.85$ )	18.36( $\pm 1.08$ )
	VAE-MAGA	11.22( $\pm 0.48$ )	178.39( $\pm 11.64$ )	18.84( $\pm 3.32$ )	213.26( $\pm 41.76$ )	5.43( $\pm 0.59$ )	8.44( $\pm 0.44$ )
	CTFG-wSP	<b>10.15</b> ( $\pm 0.56$ )	<b>133.50</b> ( $\pm 5.72$ )	<b>13.38</b> ( $\pm 1.32$ )	<b>165.96</b> ( $\pm 7.74$ )	<b>4.38</b> ( $\pm 0.14$ )	<b>7.64</b> ( $\pm 0.01$ )
	CTFG-SP	<b>7.24</b> ( $\pm 0.94$ )	<b>135.70</b> ( $\pm 16.48$ )	<b>10.23</b> ( $\pm 0.67$ )	<b>108.44</b> ( $\pm 5.82$ )	<b>2.92</b> ( $\pm 0.03$ )	<b>4.16</b> ( $\pm 0.14$ )
	CTFG-GT	<b>8.56</b> ( $\pm 0.40$ )	<b>123.02</b> ( $\pm 10.84$ )	<b>9.29</b> ( $\pm 0.34$ )	<b>114.89</b> ( $\pm 11.80$ )	<b>2.56</b> ( $\pm 0.19$ )	<b>5.26</b> ( $\pm 0.25$ )



(a) Visualization of generated images of the worst 5 samples of the R2R task. Each 1<sup>st</sup> and 2<sup>nd</sup> row of models shows the group truth samples and the generated results, respectively. The red box indicates the negative results, which do not contain all the semantics of the ground truth. We utilize randomly selected pivot images as introduced in Hwang et al. (2023) for the CTFG model.

(b) Visualization of DCI metric of the 3D Shapes dataset (training set). A more sparse matrix implies a better disentangled representation. The x-axis refers to the index of the latent vector. The y-axis represents the factor of the dataset, from 1 to 6, corresponding to floor hue, wall hue, object hue, scale, shape, and orientation.

Figure 5: Qualitative results of compositional generalization of 3D Shapes dataset.



(a) Dataset Composition (b) Consistency of Transformations (c) Factor alignment over latent vector dimension

Figure 6: Alignment between factors and latent vector dimensions. As illustrated in (a), let the coordinates of the black points be denoted as  $(x_i, y_0)$  and those of the gray points as  $(x_0, y_i)$ . The blue point then corresponds to the coordinate  $(x_0, y_0)$ , representing the intersection of independently observed components. 1) This setup captures the core intuition behind compositional generalization: novel factors can emerge at the intersection of separately observed attributes. We set black and gray points as  $(F_i, F'_0)$  and  $(F_0, F'_i)$ , respectively. Then we extend this core intuition to the latent vector space to demonstrate how each model follows this intuition and visualize it through (b). 2) Since all points are defined by their  $(x, y)$  coordinates, points with the same  $x$ -coordinate align along the  $x$ -axis, and those with the same  $y$ -coordinate align along the  $y$ -axis. Extending this idea to the latent vector space, if two samples share the same factor of variation, their latent vectors are expected to align along the axis corresponding to that factor. We visualize this core intuition in (c) to assess how well each model adheres to this structural alignment.

datasets, and the weakly supervised approach, CTFG-wSP, achieves performance close to the CTFG-GT and CTFG-SP on the challenging R2R task and across complex datasets, as shown in Table 4.

**Qualitative Results.** As illustrated in Figure 5a, both proposed methods preserve the semantics of the ground truth while baselines struggle to retain the semantics of unseen samples (ground truth). Comparing the VAE-MAGA and our models, forcing the consistency of the transformation method has a much greater impact on generalization. Additionally, our models exhibit a disentangled representation compared to the baselines, as shown in Figure 5b. This implies that a disentangled representation incorporating the symmetry structure promotes compositional generalization. Details of other results on the dSprites and MPI3D dataset are provided in Appendix E.

**Alignment of Unseen Samples.** First, as shown in Figure 6b, the baseline model yields a latent vector (red dot) that is substantially distant from the expected intersection (blue dot). In contrast, as shown in Figure 6b, the latent vector produced by our model (CTFG) is positioned much closer to the intersection, clearly outperforming the baseline. These findings suggest that our method learns a more factorized latent representation, where each dimension more effectively captures a specific factor of variation. As a result, novel combinations can be composed more meaningfully within the latent space. Second, As shown in Figure 6c,

the baseline fails to produce latent axes that are aligned across samples sharing the same factor (same color dashed lines are not parallel), and the axes do not exhibit orthogonality with respect to differing factors (different color dashed lines are not orthogonal). In contrast, our model yields latent vectors whose axes corresponding to shared factors are more consistently aligned, exhibiting a more structured and factorized representation compared to the baseline, as shown in Figure 6c.

**Stability of Generalization.** As illustrated in Figure 7, the baselines often show unstable generalization behavior during training on dSprites and 3D Shapes, whereas our model exhibits more stable improvement over training. These observations are consistent with the view that inductive bias is helpful for generalization in this setting, and further suggest that representational consistency of identical transformations contributes to improve generalization behavior.

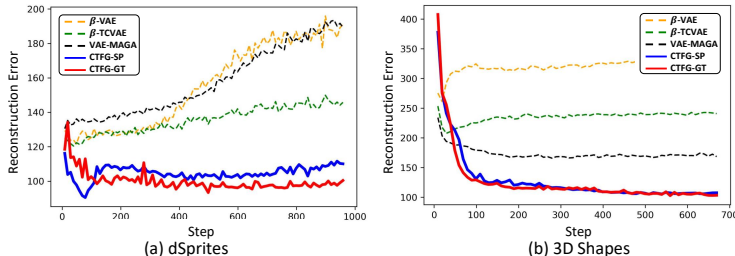


Figure 7: Test error during training.

## 7 Related Works

**Compositional Generalization.** Recent research in compositional generalization shows that models trained on disentanglement learning and verified through ground truth experimentally demonstrate that high disentangled representation does not necessarily imply compositional generalization (Montero et al., 2021; 2022; Schott et al., 2022). Differently, we consider the symmetry-based disentangled representations, recently studied in the disentanglement learning field to preserve the symmetry structure of the dataset in latent vector space. MAGANet (Hwang et al., 2023) dramatically improved compositional generalization performance by learning symmetries with the Glow model. In contrast, we study how a VAE-based model can be designed to promote more consistent representations of identical transformations.

**Disentanglement Learning.** The initially proposed methods, such as Chen et al. (2018); Kim & Mnih (2018); Higgins et al. (2017), partition each dimension to ensure mutual exclusivity, employing measures like mutual information or total correlation. However, these approaches do not account for the symmetry structure of the dataset space. Defining the symmetry group and group action as a General Linear group and matrix multiplication (Zhu et al., 2021; Jung et al., 2024; Marchetti et al., 2023) enhances disentanglement performance. Nevertheless, we theoretically demonstrate the limitations of the General Linear group for cyclic semantics in the disentangled space. Other works commonly define the symmetry group acting on the latent vector space as a cyclic group with surjective functions (Yang et al., 2021; Keurti et al., 2023; Falorsi et al., 2018). Differently, our focus is on employing isomorphism to represent the cyclic group rather than a homomorphism.

## 8 Conclusion

In this paper, we study the problem of representing identical semantic changes more consistently for disentangled representations. We analyze three commonly used group settings that can make such consistency difficult to maintain in disentangled latent spaces, and propose a Cayley-transform-based cyclic representation together with a practical learning framework. In controlled benchmark settings, the weakly-supervised variant of our design principle achieves strong disentanglement and compositional generalization, and the supervised variants further validate the potential impact of the principle. These results suggest that representational consistency of identical transformations is a useful perspective for designing inductive biases for disentanglement learning and compositional generalization simultaneously. We believe this study offers a meaningful step toward understanding how representational consistency of identical transformations can serve as a principled inductive bias, with broader implications for disentanglement learning and compositional generalization.

## References

- Nikita Balabin, Daria Voronkova, Ilya Trofimov, Evgeny Burnaev, and Serguei Barannikov. Disentanglement learning via topology. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 2474–2504. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/balabin24a.html>.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: a review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, August 2013. ISSN 0162-8828. doi: 10.1109/tpami.2013.50. URL <https://doi.org/10.1109/TPAMI.2013.50>.
- Chris Burgess and Hyunjik Kim. 3d shapes dataset. <https://github.com/deepmind/3dshapes-dataset/>, 2018.
- Jaehoon Cha and Jeyan Thiyaalingam. Orthogonality-enforced latent space in autoencoders: An approach to learning disentangled representations. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 3913–3948. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/cha23b.html>.
- Ricky T. Q. Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/1ee3dfcd8a0645a25a35977997223d22-Paper.pdf>.
- Cian Eastwood and Christopher K. I. Williams. A framework for the quantitative evaluation of disentangled representations. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=By-7dz-AZ>.
- Luca Falorsi, Pim de Haan, Tim R. Davidson, Nicola De Cao, Maurice Weiler, Patrick Forré, and Taco Cohen. Explorations in homeomorphic variational auto-encoding. *ArXiv*, abs/1807.04689, 2018. URL <https://api.semanticscholar.org/CorpusID:49672064>.
- Redre Gentner. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2):155–170, 1983. ISSN 0364-0213. doi: [https://doi.org/10.1016/S0364-0213\(83\)80009-3](https://doi.org/10.1016/S0364-0213(83)80009-3). URL <https://www.sciencedirect.com/science/article/pii/S0364021383800093>.
- Muhammad Waleed Gondal, Manuel Wüthrich, Đorđe Miladinovic, Francesco Locatello, Martin Breidt, Valentin Volchkov, Joel Bessekon Akpo, Olivier Bachem, Bernhard Scholkopf, and Stefan Bauer. On the transfer of inductive bias from simulation to the real world: a new disentanglement dataset. In *Neural Information Processing Systems*, 2019. URL <https://api.semanticscholar.org/CorpusID:182952649>.
- Irina Higgins, Loïc Matthey, Arka Pal, Christopher P. Burgess, Xavier Glorot, Matthew M. Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017.
- Irina Higgins, David Amos, David Pfau, Sébastien Racanière, Loïc Matthey, Danilo J. Rezende, and Alexander Lerchner. Towards a definition of disentangled representations. *CoRR*, abs/1812.02230, 2018. URL <http://arxiv.org/abs/1812.02230>.
- Kyle Hsu, Will Dorrell, James C. R. Whittington, Jiajun Wu, and Chelsea Finn. Disentanglement via latent quantization, 2023.
- Geonho Hwang, Jaewoong Choi, Hyunsoo Cho, and Myungjoo Kang. MAGANet: Achieving combinatorial generalization by modeling a group action. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 14237–14248. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/hwang23b.html>.

- Hee-Jun Jung, Jaehyoung Jeong, and Kangil Kim. CFASL: Composite factor-aligned symmetry learning for disentanglement in variational autoencoder. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=mDGvrH71ju>.
- Hee-Jun Jung, Jaehyoung Jeong, and Kangil Kim. Multiple invertible and partial-equivariant function for latent vector transformation to enhance disentanglement in VAEs. In *The 29th International Conference on Artificial Intelligence and Statistics*, 2026. URL <https://openreview.net/forum?id=IYZhiql7Z>.
- Hamza Keurti, Hsiao-Ru Pan, Michel Besserve, Benjamin F Grewe, and Bernhard Schölkopf. Homomorphism AutoEncoder – learning group structured representations from observed transitions. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 16190–16215. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/keurti23a.html>.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2649–2658. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/kim18b.html>.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2013. URL <https://arxiv.org/abs/1312.6114>.
- Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL [https://proceedings.neurips.cc/paper\\_files/paper/2018/file/d139db6a236200b21cc7f752979132d0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/d139db6a236200b21cc7f752979132d0-Paper.pdf).
- Erwin Kreyszig, Herbert Kreyszig, and E. J. Norminton. *Advanced Engineering Mathematics*. Wiley, Hoboken, NJ, tenth edition, 2011. ISBN 0470458364.
- Abhishek Kumar, Prasanna Sattigeri, and Avinash Balakrishnan. VARIATIONAL INFERENCE OF DISENTANGLED LATENT CONCEPTS FROM UNLABELED OBSERVATIONS. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=H1kG7GZAW>.
- Anna Kuzina, Kumar Pratik, Fabio Valerio Massoli, and Arash Behboodi. Equivariant priors for compressed sensing with unknown orientation, 2022.
- Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 4114–4124. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/locatello19a.html>.
- Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Scholkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. In *International Conference on Machine Learning*, 2020. URL <https://api.semanticscholar.org/CorpusID:211066424>.
- Giovanni Luca Marchetti, Gustaf Tegnér, Anastasiia Varava, and Danica Kragic. Equivariant representation learning via class-pose decomposition. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (eds.), *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp. 4745–4756. PMLR, 25–27 Apr 2023. URL <https://proceedings.mlr.press/v206/marchetti23b.html>.
- G. F. Marcus, S. Vijayan, S. Bandi Rao, and P. M. Vishton. Rule learning by seven-month-old infants. *Science*, 283(5398):77–80, 1999. doi: 10.1126/science.283.5398.77. URL <https://www.science.org/doi/abs/10.1126/science.283.5398.77>.

- Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.
- Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantization: VQ-VAE made simple. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=8ishA3LxN8>.
- Takeru Miyato, Masanori Koyama, and Kenji Fukumizu. Unsupervised learning of equivariant structure from sequences. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=7b7iGkuVq1Z>.
- Milton Montero, Jeffrey Bowers, Rui Ponte Costa, Casimir Ludwig, and Gaurav Malhotra. Lost in latent space: Examining failures of disentangled models at combinatorial generalisation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 10136–10149. Curran Associates, Inc., 2022.
- Milton Llera Montero, Casimir JH Ludwig, Rui Ponte Costa, Gaurav Malhotra, and Jeffrey Bowers. The role of disentanglement in generalisation. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=qbH974jKUVy>.
- Lukas Schott, Julius Von Kügelgen, Frederik Träuble, Peter Vincent Gehler, Chris Russell, Matthias Bethge, Bernhard Schölkopf, Francesco Locatello, and Wieland Brendel. Visual representation learning does not generalize strongly within the same domain. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=9RUHP1ladgh>.
- Loek Tonnaer, Luis Armando Perez Rey, Vlado Menkovski, Mike Holenderski, and Jim Portegies. Quantifying and learning linear symmetry-based disentanglement. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 21584–21608. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/tonnaer22a.html>.
- Xin Wang, Hong Chen, Yuwei Zhou, Jianxin Ma, and Wenwu Zhu. Disentangled representation learning for recommendation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):408–424, 2023. doi: 10.1109/TPAMI.2022.3153112.
- Robin Winter, Marco Bertolini, Tuan Le, Frank Noe, and Djork-Arné Clevert. Unsupervised learning of group invariant and equivariant representations. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=471pv23LDPr>.
- Tao Yang, Xuanchi Ren, Yuwang Wang, Wenjun Zeng, Nanning Zheng, and Pengju Ren. Groupifyvae: from group-based definition to vae-based unsupervised representation disentanglement. *CoRR*, abs/2102.10303, 2021. URL <https://arxiv.org/abs/2102.10303>.
- Xinqi Zhu, Chang Xu, and Dacheng Tao. Commutative lie group VAE for disentanglement learning. *CoRR*, abs/2106.03375, 2021. URL <https://arxiv.org/abs/2106.03375>.

Table 5: Notation Table

Set					
$F$	Set of latent factors	$X$	Dataset	$Z$	Set of latent vectors
$\mathcal{F}$	Space of latent factors	$\mathcal{X}$	Space of datasets	$\mathcal{Z}$	Space of latent vectors
$\mathcal{Y}$	Space of latent vectors	$\Theta$	Set of angles: $\{\theta \mid -\pi < \theta \leq \pi\}$		
Group					
$G$	Group	$G^F$	Group acted on the set of latent factors	$G^z$	Group acted on the set of latent vectors
$G^c$	Cyclic group	$g$	Group element of $G$	$\mathfrak{g}$	Lie algebra of $GL(n, \mathbb{R})$
$\alpha(\cdot, \cdot)$	Group action	$\circ_F$	Binary operation of $G_F$	$\circ_z$	Binary operation of $G_z$
$GL(n)$	General Linear group	$GL'(n)$	General Linear group implemented by matrix exponential $GL'(n) \subset GL(n)$	$g_{z_j}^{(i) \rightarrow (i+1)}$	Symmetry between $z_j^i$ and $z_j^{i+1}$
$S^1$	Circle group: $\{z \in \mathbb{C} :  z  = 1\}$	$\mathbb{R}/\mathbb{Z}$	$\{r + \mathbb{Z} \mid r \in \mathbb{R}\}$	$G \cong H$	$G$ and $H$ are isomorphic
Function					
$\Omega$	$\mathcal{F} \rightarrow \mathcal{X}$	$q_\phi$	$\mathcal{X} \rightarrow \mathcal{Z}$	$\Gamma$	$G_F \rightarrow G_z$
$\Gamma'$	$G_z \rightarrow G_y$	$e^x$	Matrix exponential	$f$	$\mathbb{R} \rightarrow \{c \mid c \in \mathbb{C},  c  = 1\}$
$f'$	$\{c \mid c \in \mathbb{C},  c  = 1\} \rightarrow \Theta$	$C.E.(.)$	Cross-entropy loss	$D_{\text{KL}}(\cdot \parallel \cdot)$	Kullback-Leibler divergence
$\mathcal{N}(\mu, \Sigma)$	Gaussian distribution	$\mathcal{U}\{a, b\}$	Discrete uniform distribution		
Linear Algebra					
$I$	Identity matrix	$V$	Codebook	$J$	Jordan normal form
$N$	Nilpotent matrix	$M^k$	Zeros matrix except $k^{\text{th}}$ column vector	$m^k$	$k^{\text{th}}$ column vector of $M^k$
$z$	Latent vector $\in \mathbb{R}^n$	$z^k$	$k^{\text{th}}$ dimension value of $z$	$c$	Latent vector $\in \mathbb{C}^n$
$c^k$	$k^{\text{th}}$ dimension value of $c$	$\theta$	$\theta \in \Theta^n$	$\theta^k$	$k^{\text{th}}$ dimension value of $\theta$
$\theta_c$	Canonical point	$\vec{0}$	Zeros vector	$\Delta v$	$\sum_k \Delta v^k$
$\Delta v^k$	Sparse vector ( $k^{\text{th}}$ dim. value $\in \mathbb{R} \setminus \{0\}$ )	$k$	Step	$k_{i \rightarrow j}$	Step from $\theta_i$ to $\theta_j$
Others					
$\mathbb{R}$	Real number	$\mathbb{Z}^+$	Integer	$\mathbb{C}$	Complex number
$\Re$	Real part of complex number	$\Im$	Imaginary part of complex number	$ F_i $	Number of factors
$\ \cdot\ $	L1 norm	$\ \cdot\ _2$	L2 norm	$l$	Ground truth

## A Details of Preliminaries

**Normal Subgroup:** A subgroup  $N$  of a group  $G$  is termed a normal subgroup of  $G$  if and only if  $gng^{-1} \in N$  for all  $g \in G$  and  $n \in N$ . Equivalently, a normal subgroup is one that is invariant under conjugation.

**Cosets:** Given a subset  $H$  of a group  $G$ , the left coset of  $H$  with respect to an element  $g \in G$  defined as:  $gH = \{gh : h \in H\} \forall g \in G$ , and the right coset  $Hg = \{hg : h \in H\} \forall g \in G$ .

**Factor (Quotient) Group:** If  $N$  is a normal subgroup of group  $G$ , the factor group (or quotient group)  $G/N$  to be the set of all left cosets of  $N$  in  $G$ :  $G/N = \{aN : a \in G\}$ .

## B Proof of Limitations

We consider the dimension-wise disentangled representation, so we constraint a few conditions as follows:

**Condition B.1.** There exists an equivariant function  $q_\phi \circ \Omega : \mathcal{F} \rightarrow \mathcal{Z}$  mapping fully disentangled factor and latent space.

**Condition B.2.**  $\mathcal{Z}$  is a  $G_z$ -set that is a symmetry group acting on  $Z$ .

**Condition B.3.** Group element  $g_z$  only affects to a single dimension value of latent vector  $z$ , where  $g_z \in G_z$ .

### B.1 Summary

**Case 1: A Limitation of  $GL(n)$**  General Linear group  $GL(n)$ , used in prior works (Zhu et al., 2021; Kuzina et al., 2022; Miyato et al., 2022; Marchetti et al., 2023), is limited in representing the representational consistency of identical transformations for disentangled representation. We first show the property of disentangled representation with matrix exponential. Then we show the limitation of  $GL(n)$ .

**Proposition B.4.** Let the symmetry group  $G_z$  ( $GL'(n)$ ) is defined as a subgroup of the General Linear group that implemented with matrix exponential, where  $GL'(n) = \{e^{\mathbf{M}} \mid \mathbf{M} \in \mathbb{R}^{n \times n}\}$ ,  $g^k$  is an element of  $GL'(n)$ , and  $g = \prod_k g^k$ . Then  $e^g z = e^{\mathbf{I}g} z + v'$ .

**Theorem B.5.** (Limit of  $GL'(n)$ ) According to Proposition Theorem B.4, only the identity matrix ( $g = \mathbf{I}$ ) represents the cyclic group of the dataset with representational consistency of identical transformations, where  $g \in GL'(n)$ .

**Theorem B.6.** (*Limit of  $GL(n)$* ) If  $H \subset GL(n)$ , then only the identity matrix of  $GL(n)$  represents the cyclic group of the dataset with representational consistency of identical transformations, where  $H = \{h|h = \mathbf{I} + \mathbf{M}^k\}$ ,  $\mathbf{m}^k$  is a column vector of  $\mathbf{M}^k$ ,  $\mathbf{m}^j = \mathbf{0}$  and  $j \in \{1, 2, \dots, n\} \setminus \{k\}$ .

Therefore, the limitation of  $GL(n)$  is that only the  $\mathbf{I}$  represents the representational consistency of identical transformations with disentangled representation according to the Theorem Theorem B.5, and Theorem B.6. It implies that if  $g \neq \mathbf{I}$ , then  $g$  can not represent the representational consistency of identical transformations. More details are in Appendix Section B.2.

**Case 2: A Limitation of Using Vector Addition** Another setting that causes problems in maintaining representational consistency of identical transformations for cyclic group using vector addition is utilized to define a group action between two latent vectors for an equivariant model, as used in Hwang et al. (2023); Balabin et al. (2024).

**Corollary B.7.** If the group action is defined as  $\alpha(g, \mathbf{z}_i) = g + \mathbf{z}_i$ , then only zero vector represent the representational consistency of identical transformations for cyclic group with disentangled representation, where  $\mathbf{z} \in \mathbb{R}^n$ .

According to the Theorem Theorem B.7, it also shows a limitation in that only the identity element  $\mathbf{0}$  represents the consistent symmetry. More details of the proof are in Appendix Section B.3.

**Case 3: A Limitation of Surjective Function** The last setting is a surjective function that maps latent vectors to the unit circle (Yang et al., 2021; Tonnaer et al., 2022; Cha & Thiyagalingam, 2023), causing undifferentiated symmetries under more general latent factors of variation and losing part of the dataset’s symmetry information.

**Corollary B.8.** By the equivariant and surjective function  $b: \mathcal{Z} \rightarrow \mathcal{Y}$ , the capacity of  $\mathcal{Z}$  and  $\mathcal{Y}$  is  $|\mathcal{Z}| \geq |\mathcal{Y}|$  then  $\Gamma'$  is an endomorphism because  $|G_{\mathcal{Z}}| \geq |G_{\mathcal{Y}}|$ . On the other hand, isomorphism identically maps the two spaces ( $|G_{\mathcal{Z}}| = |G_{\mathcal{Y}}|$ ).

Therefore, if  $b$  is surjective and not injective, then there exists at least one case where  $\Gamma'(g_i) = \Gamma'(g_j)$ . It implies that loss of symmetry structure occurs with a surjective function. More details of the proof are in Appendix Section B.4.

## B.2 Proof: Limitations of $GL(n)$

**Proof: Limitations of  $GL(n)$ , implemented by the matrix exponential, to Represent the Cyclic Group on the Disentangled Space**

**Condition B.9.** The symmetry group  $G_z$  ( $GL'(n)$ ) acting on latent vector space is defined as a subgroup of the General Linear group, implemented with matrix exponential.

**Condition B.10.** For  $g^k \in \mathbb{R}^{n \times n}$  and  $g = \prod_k g^k$ ,  $g^k$  only affects the  $k^{th}$  dimension value of vector  $\mathbf{z}$ .

In prior works (Jung et al., 2024; Zhu et al., 2021; Kuzina et al., 2022; Miyato et al., 2022; Marchetti et al., 2023), the General Linear group  $GL(n)$  is usually implemented with the Lie algebra  $\mathfrak{g}$  to represent the symmetries between two inputs in the latent vector space:

$$g = e^{\sum_i \alpha_i \mathfrak{g}_i}, \quad (11)$$

where  $g \in GL(n)$ ,  $\alpha \in \mathbb{R}$ ,  $\mathfrak{g} \in \mathbb{R}^{n \times n}$ , and the matrix exponential  $e^{\mathfrak{g}}$  defined as  $e^{\mathfrak{g}} = \sum_{n=0}^{\infty} \frac{1}{n!} \mathfrak{g}^n$ . In addition, group  $GL(n)$  acts on the latent vector space  $\mathcal{Z}$  with group action:

$$\alpha(g, \mathbf{z}) = g\mathbf{z}, \quad (12)$$

where latent vector  $\mathbf{z} \in \mathbb{R}^n$ , commonly used in previous works (Zhu et al., 2021; Jung et al., 2024; Kuzina et al., 2022; Marchetti et al., 2023). We first show the property of disentangled representation with matrix exponential.

**Proposition B.11.** *Let the symmetry group  $G_z$  ( $GL'(n)$ ) is defined as a subgroup of the General Linear group that implemented with matrix exponential, where  $GL'(n) = \{e^{\mathbf{M}} | \mathbf{M} \in \mathbb{R}^{n \times n}\}$ ,  $g^k$  is an element of  $GL'(n)$ , and  $g = \prod_k g^k$ . Then  $e^g \mathbf{z} = e \mathbf{I} g \mathbf{z} + \mathbf{v}'$ .*

*Proof.* If group element  $g$  acts on latent vector then,  $g \mathbf{z} - \mathbf{z} = \Delta \mathbf{v}$ , where  $\Delta \mathbf{v} = \sum_k \Delta \mathbf{v}^k$ ,  $\Delta \mathbf{v}^k$  is a sparse vector ( $k^{\text{th}}$  dimension value  $\in \mathbb{R} \setminus \{0\}$ , otherwise it is zero), and  $g^k \mathbf{z} - \mathbf{z} = \Delta \mathbf{v}^k$ . Then we define  $(g)^n \mathbf{z} - \mathbf{z} = n \Delta \mathbf{v} + (n-1) \mathbf{v}'_n$ , where  $\mathbf{v}'_n \in \mathbb{R}^n$  is an arbitrary real vector. Group element  $g$  represents the cyclic semantics of the dataset space, then it satisfies the following equation:

$$\begin{aligned} (g_i - \mathbf{I}) \mathbf{z} &= \Delta \mathbf{v} \\ \frac{1}{2!} ((g_i)^2 - \mathbf{I}) \mathbf{z} &= \frac{2}{2!} (\Delta \mathbf{v} + \frac{1}{2} \mathbf{v}'_2) \\ &\vdots \\ \lim_{n \rightarrow \infty} \frac{1}{n!} ((g_i)^n - \mathbf{I}) \mathbf{z} &= \lim_{n \rightarrow \infty} \frac{1}{(n-1)!} (\Delta \mathbf{v} + \frac{1}{n} \mathbf{v}'_n). \end{aligned} \quad (13)$$

By adding left- and right-hand side of Eq. Equation (13), we then get:

$$\begin{aligned} &\Rightarrow (e^{g_i} - \mathbf{I}) \mathbf{z} - (e-1) \mathbf{I} \mathbf{z} = e \mathbf{I} \Delta \mathbf{v} + \mathbf{v}' \\ &\Rightarrow e^{g_i} \mathbf{z} = e \mathbf{I} g_i \mathbf{z} + \mathbf{v}', \end{aligned} \quad (14)$$

where  $g_i \in G$ ,  $\mathbf{v}' = \lim_n \sum_n \frac{1}{n!} \mathbf{v}'_n$  and  $\mathbf{v}' = \mathbf{v}' + e \mathbf{I} \Delta \mathbf{v}$ .  $\square$

**Lemma B.12.** *By the Proposition Theorem B.11, if  $\mathbf{v}' = \vec{0}$  in Eq. Equation (14), then  $g = \mathbf{I}$  and the index of the nilpotent matrix of Jordan normal form of  $\mathbf{g}$  is 2.*

*Proof.* The trivial solution of  $(e^g - e \mathbf{I} g) \mathbf{z} = 0, \forall \mathbf{z} \in \mathcal{Z}$  is that

$$e^g - e \mathbf{I} g = 0. \quad (15)$$

Every matrix  $\mathbf{g} \in \mathbb{C}^{n \times n}$  has a Jordan normal form  $\mathbf{J}$  as  $\mathbf{g} = \mathbf{S} \mathbf{J} \mathbf{S}^{-1}$ . Then group element ( $g = e^{\mathbf{g}}$ ) follows as:

$$\begin{aligned} e^{\mathbf{g}} &= \lim_{n \rightarrow \infty} \mathbf{I} + \mathbf{S} \mathbf{J} \mathbf{S}^{-1} + \frac{1}{2!} \mathbf{S} \mathbf{J}^2 \mathbf{S}^{-1} + \dots + \frac{1}{n!} \mathbf{S} \mathbf{J}^n \mathbf{S}^{-1} \\ &= \lim_{n \rightarrow \infty} \mathbf{I} + \mathbf{S} (\mathbf{J} + \frac{1}{2!} \mathbf{J}^2 + \dots + \frac{1}{n!} \mathbf{J}^n) \mathbf{S}^{-1} \\ &= \mathbf{I} + \mathbf{S} (e^{\mathbf{J}} - \mathbf{I}) \mathbf{S}^{-1} \\ &= \mathbf{S} e^{\mathbf{J}} \mathbf{S}^{-1} \end{aligned} \quad (16)$$

In the same way, the exponential of  $g$  is equal to:

$$e^g = \mathbf{S} e^{e^{\mathbf{J}}} \mathbf{S}^{-1} \quad (17)$$

Therefore group element  $g$  satisfies

$$\begin{aligned} \mathbf{S} e^{e^{\mathbf{J}}} \mathbf{S}^{-1} &= e \mathbf{I} \mathbf{S} e^{\mathbf{J}} \mathbf{S}^{-1} \\ &\Rightarrow e^{e^{\mathbf{J}}} = e^{\mathbf{I}} e^{\mathbf{J}} \quad (\because e^{\mathbf{I}} = e^{\mathbf{I}}) \\ &\Rightarrow e^{e^{\mathbf{J}}} = e^{\mathbf{I} + \mathbf{J}} \quad (\because \mathbf{I} \mathbf{J} = \mathbf{J} \mathbf{I}) \\ &\therefore e^{\mathbf{J}} = \mathbf{I} + \mathbf{J} \end{aligned} \quad (18)$$

If  $\mathbf{J} = \mathbf{0}$ , then  $g$  satisfies the Eq. Equation (18) and  $g = \mathbf{I}$ . If  $\mathbf{J} \neq \mathbf{0}$  then,

$$e^{\mathbf{J}} = \begin{bmatrix} e^{\lambda_1} & e_{1,2}^{\mathbf{J}} & \cdots & e_{1,n}^{\mathbf{J}} \\ & e^{\lambda_2} & \cdots & e_{2,n}^{\mathbf{J}} \\ & & \ddots & \vdots \\ & & & e^{\lambda_n} \end{bmatrix}, \text{ and } \mathbf{I} + \mathbf{J} = \begin{bmatrix} \lambda_1 + 1 & (\mathbf{I} + \mathbf{J})_{1,2} & \cdots & (\mathbf{I} + \mathbf{J})_{1,n} \\ & \lambda_2 + 1 & \cdots & (\mathbf{I} + \mathbf{J})_{2,n} \\ & & \ddots & \vdots \\ & & & \lambda_n + 1 \end{bmatrix}, \quad (19)$$

where empty values are all zero. To satisfy the Eq. Equation (18),  $\lambda_i = 0$  for  $e_i^\lambda = \lambda_i + 1$ , then  $\mathbf{J} = \mathbf{D} + \mathbf{N} = \mathbf{N}$  because diagonal of  $\mathbf{D}$   $\lambda_i = 0$ , where  $\mathbf{D}$  is a diagonal matrix and  $\mathbf{N}$  is a nilpotent matrix. Therefore,

$$\begin{aligned} e^{\mathbf{J}} &= e^{\mathbf{N}} \text{ and } \mathbf{I} + \mathbf{J} = \mathbf{I} + \mathbf{N} \\ &\Rightarrow e^{\mathbf{N}} = \mathbf{I} + \mathbf{N} \\ &\Rightarrow \lim_{n \rightarrow \infty} \frac{1}{2!} \mathbf{N}^2 + \frac{1}{3!} \mathbf{N}^3 + \dots + \frac{1}{n!} \mathbf{N}^n = 0 \end{aligned} \quad (20)$$

Therefore if the index of nilpotent matrix is 2 and  $e_{i,j}^{\mathbf{J}} = (\mathbf{I} + \mathbf{J})_{i,j}$  where  $i < j$ , then it satisfies the Eq. Equation (15) □

**Lemma B.13.** *If the index of the nilpotent matrix of Jordan normal form of  $\mathfrak{g}$  is 2, then  $g = e^{\mathfrak{g}}$  does not represent the element of cyclic group.*

*Proof.* If  $g$  is an element of cyclic group then there exists  $n^{\text{th}}$  power of  $g$  such that  $g^n$  is the identity matrix.

$$\begin{aligned} g^n &= \mathbf{S}(\mathbf{I} + n\mathbf{N})\mathbf{S}^{-1} = \mathbf{I} \quad (\because \mathbf{N}^2 = \mathbf{0}) \\ &\Rightarrow \mathbf{I} + \mathbf{S}n\mathbf{N}\mathbf{S}^{-1} = \mathbf{I} \\ &\Rightarrow \mathbf{S}n\mathbf{N}\mathbf{S}^{-1} = \mathbf{0}. \end{aligned} \quad (21)$$

To satisfy the Eq. Equation (21),  $\mathbf{N} = \mathbf{0}$  because the index of  $\mathbf{N}$  is 2. □

By the Lemma Theorem B.12 and Theorem B.13, there exists only one element to represent the element of cyclic group of the dataset in the disentangled space.

**Lemma B.14.** *If  $\mathbf{v}' \in \mathbb{R}^n \setminus \{0\}$ , then 1)  $g^k$  represents the element of cyclic group as  $g^k = \mathbf{I} + \mathbf{M}^k$  and  $\mathbf{m}^k \in \mathbb{R}^n$ , where  $\mathbf{m}^k$  is a  $k^{\text{th}}$  column vector of  $\mathbf{M}^k$ ,  $\mathbf{m}^j = \mathbf{0}$  and  $j \in \{1, 2, \dots, n\} \setminus \{k\}$ .*

*Proof.* As we define that  $g^k$  only affect to a single dimension value of  $\mathbf{z}$  we rewrite Eq. Equation (14) as follows:

$$e^{g^k} \mathbf{z} = e\mathbf{I}g^k \mathbf{z} + \mathbf{v}'^k, \quad (22)$$

where  $\mathbf{v}'^k$  is a sparse vector. For Eq. Equation (22)  $\forall \mathbf{z} \in \mathcal{Z}$  then symmetry  $g^k$  follows as:

$$e^{g^k} - e\mathbf{I}g^k = [\vec{0} \quad \dots \quad \vec{0} \quad \mathbf{m}^k \quad \vec{0} \quad \dots \quad \vec{0}], \quad (23)$$

where  $\vec{0}$  is a zero column vector.

Then, satisfying the Eq. Equation (23) and affecting a single dimension:

$$\begin{aligned} \mathbf{z}^\top g^k &= [z_1, \dots, z_{k-1}, z_k + \alpha, z_{k+1}, \dots, z_n] \\ \mathbf{z}^\top e\mathbf{I}g^k &= [ez_1, \dots, ez_{k-1}, e(z_k + \alpha), ez_{k+1}, \dots, ez_n] \\ \therefore \mathbf{z}^\top e^{g^k} &= [ez_1, \dots, ez_{k-1}, z_k + \beta, ez_{k+1}, \dots, ez_n] \\ &\Rightarrow \mathbf{z}^\top (e^{g^k} - e\mathbf{I}) = [0, \dots, 0, (1 - e)z_k + \beta, \dots, 0]. \end{aligned} \quad (24)$$

For Eq. Equation (24) for all  $\mathbf{z}$ , then

$$\begin{aligned} e^{g^k} - e\mathbf{I} &= [\vec{0} \quad \dots \quad \vec{0} \quad \mathbf{m}^k \quad \vec{0} \quad \dots \quad \vec{0}] \\ \therefore e^{g^k} &= e\mathbf{I} + [\vec{0} \quad \dots \quad \vec{0} \quad \mathbf{m}^k \quad \vec{0} \quad \dots \quad \vec{0}]. \end{aligned} \quad (25)$$

By the same way,  $g^k = \mathbf{I} + [\vec{0} \quad \dots \quad \vec{0} \quad \mathbf{m}^k \quad \vec{0} \quad \dots \quad \vec{0}] = \mathbf{I} + \mathbf{M}^k$  because  $\mathbf{z}^\top (g^k - \mathbf{I})$  is a sparse vector. □

**Lemma B.15.** For  $g_k$  to represent element of cyclic group,  $g_k = \mathbf{I}$ .

*Proof.* If  $\mathbf{I} + \mathbf{M}^k$  represents the element of cyclic group then, there exists a  $n$  where  $(\mathbf{I} + \mathbf{M}^k)^{(n+1)} = \mathbf{I}$ . The formation of the power of the matrix  $(\mathbf{I} + \mathbf{M}^k)^n$  is also  $\mathbf{I} + \mathbf{M}_n^k$ , where  $\mathbf{M}_n^k \in \mathbf{M}^k$  is an arbitrary real matrix. Therefore  $\mathbf{M}_n^k = \mathbf{0}$ , then  $g^k = \mathbf{I}$ .  $\square$

It implies that  $g^k$  represents only an identity transformation of dataset space.

**Theorem B.16.** (*Limit of  $GL'(n)$* ) According to Proposition Theorem B.11, only the identity matrix ( $g = \mathbf{I}$ ) represents the cyclic group of the dataset with representational consistency of identical transformations, where  $g \in GL'(n)$

*Proof.* Through the Lemma Theorem B.12 to Lemma Theorem B.15, if the group element  $g$  represents the element of cyclic group, then only the identity matrix satisfies the Eq. Equation (13). There is always a case as  $g = e^{\mathfrak{g}}$ , and  $(g)^{n-1} = (e^{\mathfrak{g}})^{n-1}$  but  $(g)^n \neq (e^{\mathfrak{g}})^n$ , where  $g \neq \mathbf{I}$ , because  $g$  can not represent the element of cyclic group. By the equivariant function  $q_\phi$ :

$$\begin{aligned} q_\phi(x_1) &= z_1 \\ q_\phi(g_x \circ_x x_1) &= \Gamma(g_x) \circ_z z_1 \\ &\vdots \\ q_\phi(g_x^{n-1} \circ_x x_1) &= [\Gamma(g_x)]^{n-1} \circ_c z_1 \\ q_\phi(g_x^n \circ_x x_1) &\neq [\Gamma(g_x)]^n \circ_c z_1 (\because [\Gamma(g_x)]^n \neq (e^{\mathfrak{g}})^n). \end{aligned} \tag{26}$$

It implies that element of cyclic group  $g_x$  between the  $x_k$  and  $x_{k+1}$  is divided as:

$$\Gamma(g_x) = \begin{cases} e^{\mathfrak{g}} & \text{if } k < n \\ e^{\mathfrak{g}'} & \text{if } k = n \end{cases}, \text{ where } \mathfrak{g} \neq \mathfrak{g}'. \tag{27}$$

$\square$

Therefore, representing the cyclic group of the dataset with representational consistency of identical transformations is impossible according to the Theorem B.16.

### Limitations of $GL(n)$ to Represent the Cyclic Group on the Disentangled Space

**Theorem B.17.** (*Limit of  $GL(n)$* ) If  $H \subset GL(n)$ , then only the identity matrix of  $GL(n)$  represents the cyclic group of the dataset with representational consistency of identical transformations, where  $H = \{h | h = \mathbf{I} + \mathbf{M}^k\}$ ,  $\mathbf{m}^k$  is a column vector of  $\mathbf{M}^k$ ,  $\mathbf{m}^j = \vec{0}$  and  $j \in \{1, 2, \dots, n\} \setminus \{k\}$ .

*Proof.* If the invertible matrix  $g$  changes a single dimension value of the latent vector  $\mathbf{z}$  then,

$$\begin{aligned} \mathbf{z}^\top g - \mathbf{z} &= [0, \dots, 0, \alpha, 0, \dots, 0] \\ \Rightarrow \mathbf{z}^\top (g - \mathbf{I}) &= [0, \dots, 0, \alpha, 0, \dots, 0]. \\ \Rightarrow [z_1, \dots, z_n] \begin{bmatrix} g_{11} - 1 & \cdots & g_{1n} \\ \vdots & \ddots & \vdots \\ g_{n1} & \cdots & g_{nn} - 1 \end{bmatrix} &= [0, \dots, 0, \alpha, 0, \dots, 0] \end{aligned} \tag{28}$$

As defined the  $G$ -set as a latent vector space  $\mathcal{Z}$ , group element  $g$  satisfies the Eq. Equation (28) over all vectors. Then  $g - \mathbf{I}$  elements are all 0, except the  $k^{\text{th}}$  column vector ( $\mathbf{m}^k \neq 0$ ). Therefore, group element  $g \in H$ .  $\square$

**Theorem B.18.** If  $h \in H \setminus \{\mathbf{I}\}$ , then this set does not represent the representational consistency of identical transformations with element of cyclic group.

*Proof.* According to Lemma Theorem B.15,  $h^k = \mathbf{I}$ . Therefore,  $h^k$  represents only the representational consistency of identical transformations of the dataset space.  $\square$

Therefore, representing the cyclic group of the dataset with representational consistency of identical transformations implemented by the  $GL(n)$  is impossible except in the case of the identity matrix.

### B.3 Proof: Limitations of Vector Addition for Cyclic Group on the Disentangled Space

In the previous works (Hwang et al., 2023), the vector addition is used to define a group action between two latent vectors for an equivariant model, where the group action  $\alpha(g, \mathbf{z}) = g + \mathbf{z}$ .

**Theorem B.19.** *If the group action is defined as  $\alpha(g, \mathbf{z}_i) = g + \mathbf{z}_i$ , then  $g$  does not represent the representational consistency of identical transformations for element of cyclic group with disentangled representation, where  $g \in G \setminus \{\vec{0}\}$  and  $\mathbf{z} \in \mathbb{R}^n$ .*

*Proof.* If  $g$  represents the element of cyclic group of  $\mathcal{X}$ , then there exists:

$$(g)^n = ng = \vec{0}. \quad (29)$$

The solution of Eq. Equation (29) is  $g = \vec{0}$ . There is always a case as  $g = \mathbf{a}$ , and  $(g)^{n-1} = (n-1)\mathbf{a}$  but  $(g)^n \neq n\mathbf{a}$ , where  $g \neq \vec{0}$ , because  $g$  can not represent the element of cyclic group. By the equivariant function  $q_\phi$ :

$$\begin{aligned} q_\phi(x_1) &= z_1 \\ q_\phi(g_x \circ_x x_1) &= \Gamma(g_x) \circ_z z_1 \\ &\vdots \\ q_\phi(g_x^{n-1} \circ_x x_1) &= (n-1)[\Gamma(g_x)] \circ_c z_1 \\ q_\phi(g_x^n \circ_x x_1) &\neq n[\Gamma(g_x)] \circ_c z_1 \quad (\because n[\Gamma(g_x)] \neq n\mathbf{a}). \end{aligned} \quad (30)$$

It implies that element of cyclic group  $g_x$  between the  $x_k$  and  $x_{k+1}$  is divided as:

$$\Gamma(g_x) = \begin{cases} \mathbf{a} & \text{if } k < n \\ \mathbf{a}' & \text{if } k = n \end{cases}, \text{ where } \mathbf{a} \neq \mathbf{a}'. \quad (31)$$

$\square$

According to Theorem B.19, the group action defined as  $\alpha(g, \mathbf{z}) = g + \mathbf{z}$  represents the cyclic group of the dataset with representational consistency of identical transformations when  $g = \vec{0}$ . However, the group elements are insufficient to encompass the entire cyclic group of input space. Additionally, this causes the inconsistency issue when holding consistency with vector addition.

### B.4 Proof: Loss of Symmetry Structure with Endomorphism

**Definition B.20.** Let  $(G, \cdot), (H, \circ)$  be two groups. If mapping function  $\Gamma : G \rightarrow H$ , s.t.  $\Gamma(g_i \cdot g_j) = \Gamma(g_i) \circ \Gamma(g_j)$ , then  $\Gamma$  is called homomorphism.

**Definition B.21.** Let  $\Gamma$  is surjective then  $\Gamma$  is called endomorphism.

**Definition B.22.** Let a special case of homomorphism where  $\Gamma$  is bijective is called isomorphism.

**Corollary B.23.** *By the equivariant and surjective function  $b : \mathcal{Z} \rightarrow \mathcal{Y}$ , the capacity of  $\mathcal{Z}$  and  $\mathcal{Y}$  is  $|\mathcal{Z}| \geq |\mathcal{Y}|$  then  $\Gamma$  is an endomorphism because  $|G_{\mathcal{Z}}| \geq |G_{\mathcal{Y}}|$ . On the other hand, isomorphism identically maps the two spaces ( $|G_{\mathcal{Z}}| = |G_{\mathcal{Y}}|$ ).*

Therefore, if  $b$  is surjective and not injective, then there exists at least one case where  $\Gamma'(g_i) = \Gamma'(g_j)$ . It implies that loss of symmetry structure occurs with a surjective function.

## C Details of the Proposed Method Motivation

### C.1 Why a Bijective Cyclic Representation?

**From Factor States to Cyclic Group Structure.** As established in Section 4, the symmetry group acting on the factor states is  $G^F \cong \mathbb{Z}_{|F^1|} \times \cdots \times \mathbb{Z}_{|F^n|}$ . A cyclic group  $G = \{e, g_1, g_2, \dots, g_{n-1}\}$  has the convenient property that all elements are integer powers of a single generator; consequently, if the model learns one generator, it implicitly represents the entire group. Motivated by Yang et al. (2021), we implement the cyclic group as the  $N$ -th roots of unity. Concretely, the cyclic group is

$$G^c = G_1^c \times G_2^c \times \cdots \times G_k^c, \quad G_i^c = \{g_i^c \mid g_i^c = \frac{2\pi}{N}k, k \in \{0, 1, \dots, N-1\}\}, \quad (32)$$

where  $N \in \mathbb{Z}^+$ . The group action is defined as  $\alpha : \Theta^n \times G^c \rightarrow \Theta^n$ ,  $\alpha(g^c, \theta) = g^c + \theta$ , where  $\theta \in \Theta^n$  and  $-\pi < \theta_i^k \leq \pi$ .

**Isomorphism, not Homomorphism.** The surjective-mapping family (Case 3 in Section 4) represents cyclic symmetries via a homomorphism  $G^z \rightarrow G^y$ , which can collapse distinct elements and lose group structure. To avoid this, we require an isomorphism  $G^z \cong G^F$ , i.e., a bijective structure-preserving map between the latent group and the factor group. This isomorphism requirement directly motivates working on the unit circle  $S^1$  via a bijective (rather than surjective) embedding.

**Connecting to Definition 3.1.** The key property of the bijective cyclic representation is that identical transformations map to identical group elements regardless of the starting state. We formalize this as the following proposition.

**Proposition C.1** (Consistent Representation via Bijective Cyclic Embedding). *Let  $f' \circ f : \mathbb{R} \rightarrow \Theta$  be the bijective map defined in Section 5.2 (Cayley transform composed with angle extraction), and let  $\theta^k = \arg \min_{V^i \in \mathbf{V}} |V^i - \theta^k|$  be the nearest-neighbor projection onto the fixed grid  $V$  defined in Section 5.3. Then, under the cyclic group structure  $G^{F^i} \cong \mathbb{Z}_N$ , the resulting latent representation satisfies Definition 3.1: for any two pairs  $(f_{v_i}, f_{v_j})$  and  $(f_{v_n}, f_{v_m})$  with  $f_{v_j} = g^F \cdot f_{v_i}$  and  $f_{v_m} = g^F \cdot f_{v_n}$ , the corresponding latent transformation  $\Gamma(g^F) = \frac{2\pi}{N}k$  depends only on  $g^F$  (equivalently, on the step  $k$ ), and not on the particular pair.*

*Proof sketch.* Because  $f' \circ f$  is bijective, each real-valued latent  $z$  maps to a unique angle  $\theta \in \Theta^n$ . The fixed grid  $\mathbf{V}$  partitions  $\Theta$  into  $N$  equal-width bins of size  $\frac{2\pi}{N}$ , so any angle is snapped to one of the  $N$  canonical grid points. The group action  $\alpha(g^c, \theta) = g^c + \theta$  with  $g_i^c = \frac{2\pi}{N}k$  then advances the grid point by exactly  $k$  steps. Since the grid is fixed (not sample-dependent), the same  $k$ -step transformation always produces the same latent shift  $\frac{2\pi}{N}k$ , independently of the starting angle  $\theta_i$  or  $\theta_n$ . This directly satisfies the pair-independence condition in Definition 3.1.  $\square$

## D Details of Experiments Setting

### D.1 Resources

We set the following settings for all experiments on a single Galaxy 2080Ti GPU, a single Galaxy 3090 GPU, and a single NVIDIA A100 GPU for the dSprites 3D Shapes and MPI3D datasets. The Python version is 3.7.10, and the PyTorch version is 1.9.1. More details are in the README.md file.

### D.2 Datasets

1) The dSprites dataset consists of 737,280 binary  $64 \times 64$  images with five independent ground truth factors (number of values), i.e. x-position (32), y-position (32), orientation (40), shape (3), and scale (6) (Matthey et al., 2017). Any composite transformation of x- and y-position, orientation (2D rotation), scale, and shape is commutative. 2) The 3D Shapes dataset consists of 480,000 RGB  $64 \times 64 \times 3$  images with six independent ground truth factors: orientation (15), shape (4), floor color (10), scale (8), object color (10), and wall color (10) (Burgess & Kim, 2018). 3) The MPI3D (real-world complex) dataset consists of 460,800 RGB  $64 \times 64 \times 3$

images with seven independent ground truth factors: color (4) shape (4), size (2), height (3), background color (3), horizontal (40), and vertical axis (40) (Gondal et al., 2019). Additionally, we use cLPR dataset<sup>1</sup> consists of 250,047 RGB  $64 \times 64 \times 3$  images with three independent ground truth factors: x-rotation (63), y-rotation (63), and z-rotation (63).

### D.3 Setting for compositional Generalization

**Train and Test datasets** We except the case [shape=*ellips*, position-x  $\geq 0.6$ , position-y  $\geq 0.6$ ,  $120^\circ \leq$  rotation  $\leq 240^\circ$ , scale  $< 0.6$ ] for dSprites r2e training set and [shape=*square*, position-x  $\geq 0.5$ ] for dSprites r2r training set.

We except the case [floor-hue  $> 0.5$ , wall-hue  $> 0.5$ , object-hue  $\geq 0.5$ , shape=*cylinder*, object-scale=1, object-orientation=0] for 3D Shapes r2e training set and [object-hue  $\geq 0.5$ , shape=*oblong*] for 3D Shapes r2r training set.

We except the case [shape = *cone*, object size = 0, cameraheight = 1, background color = purple, object color  $\in$  {blue, brown, olive}, horizontal axis  $\geq 20$ , vertical axis  $\geq 20$ ] for r2e training set and [shape = *cylinder*, scale  $< 6$ , orientation,  $16 \leq$  horizontal axis  $< 32$ , vertical axis] for r2r training set.

**Hyper-Parameter Tuning** We set  $\beta \in \{1, 2, 10\}$  for  $\beta$ -VAE (Higgins et al., 2017) and  $\beta$ -TCVAE (Chen et al., 2018), and  $\beta \in \{1, 10\}$  for VAE-MAGA, which employs the MAGA-net proposed module on the CNN-based encoder and decoder, instead of Glow (Kingma & Dhariwal, 2018). We set the common hyper-parameters of the proposed method at  $\alpha \in \{100, 1000\}$ ,  $\gamma \in \{1, 10\}$ ,  $\lambda \in \{0.0, 1.0\}$  for the supervised and ground truth models, and  $\beta \in \{1.0, 2.0\}$  for the supervised method. We run each model with three seeds,  $\in \{1, 2, 3\}$ . We set  $N = 1000$  and  $n' = 10$ .

**Decoder Equivariant Loss** For compositional generalization, we add the decoder equivariant loss as:

$$\mathcal{L}_{de} = R.E(x_j, p_\theta \circ g_{i \rightarrow j} \circ_\theta (f' \circ f \circ q_\phi(x_i))), \quad (33)$$

where  $R.E(\cdot)$  is a reconstruction error and  $p_\theta$  is a decoder. We add the  $\mathcal{L}_{de}$  to the objective losses with hyper-parameter  $\lambda$ .

### D.4 Setting for Disentanglement Learning

**Hyper-Parameter Tuning** We set  $\beta \in \{1, 2, 10\}$  for  $\beta$ -VAE (Higgins et al., 2017) and  $\beta$ -TCVAE (Chen et al., 2018),  $\gamma \in \{10, 20, 40\}$  for Factor-VAE (Kim & Mnih (2018)),  $hy_{rec} \in \{0.1, 0.2, 0.7\}$  for CLG-VAE (Zhu et al., 2021),  $\beta = 1$  for Ada-GVAE (Locatello et al., 2020). We set the common hyper-parameters of proposed method  $\alpha \in \{100, 1000\}$ ,  $\gamma = 1$  for supervised and ground truth model,  $\beta \in \{1.0, 2.0\}$  for supervised method, and  $\lambda = 1.0, p = 0.5$  for weakly-supervised method. We run 10 seed variance over each model with seed  $\in \{1, 2, \dots, 10\}$ . We set  $N = 1000$  and  $n' = 10$ . We evaluate four metrics  $\beta$ -VAE metric (Higgins et al., 2017), Factor VAE metric (Kim & Mnih, 2018), SAP (Kumar et al., 2018), and DCI (Eastwood & Williams, 2018).

## E compositional Generalization Results

As shown in Figure 8a, we selected the four worst samples (those with the highest reconstruction errors): 1)  $\beta$ -VAE results contain only position semantics, 2)  $\beta$ -TCVAE captures the position and scale values but fails to capture the shape and rotation factors, 3) VAE-MAGA struggles with generalization. Even though our method does not capture all semantics, it shows improvement compared to the baselines: 4) the supervised method misses either the shape or rotation, and 5) the GT model only misses the shape semantic. As shown in Figure 8b, the representations of the baseline are not close to a disentangled representation. In contrast, the representation of the supervised method approaches a disentangled representation and shows better generalization. This implies that a disentangled representation containing the symmetry structure could benefit compositional generalization.

<sup>1</sup><https://github.com/yvan/cLPR>

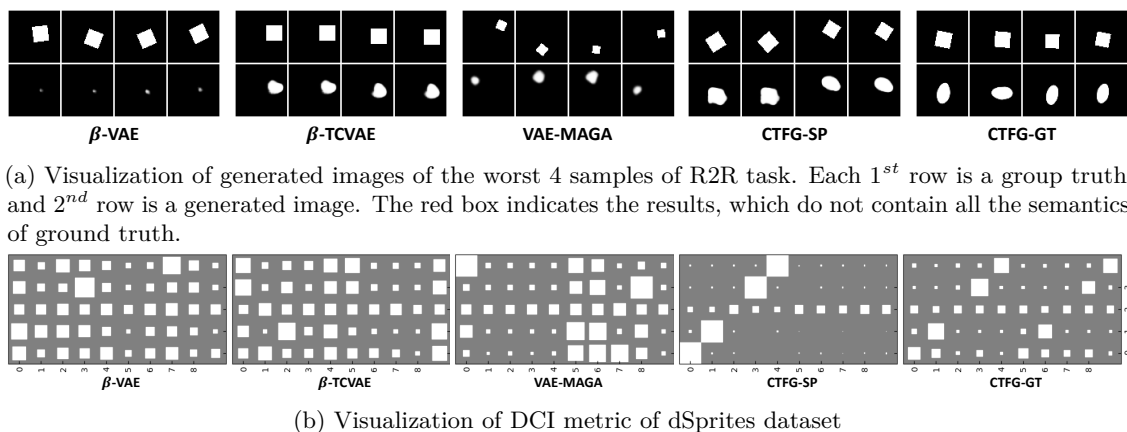


Figure 8: Qualitative results of compositional generalization of dSprites dataset (R2R). A more sparse matrix implies clear disentanglement.

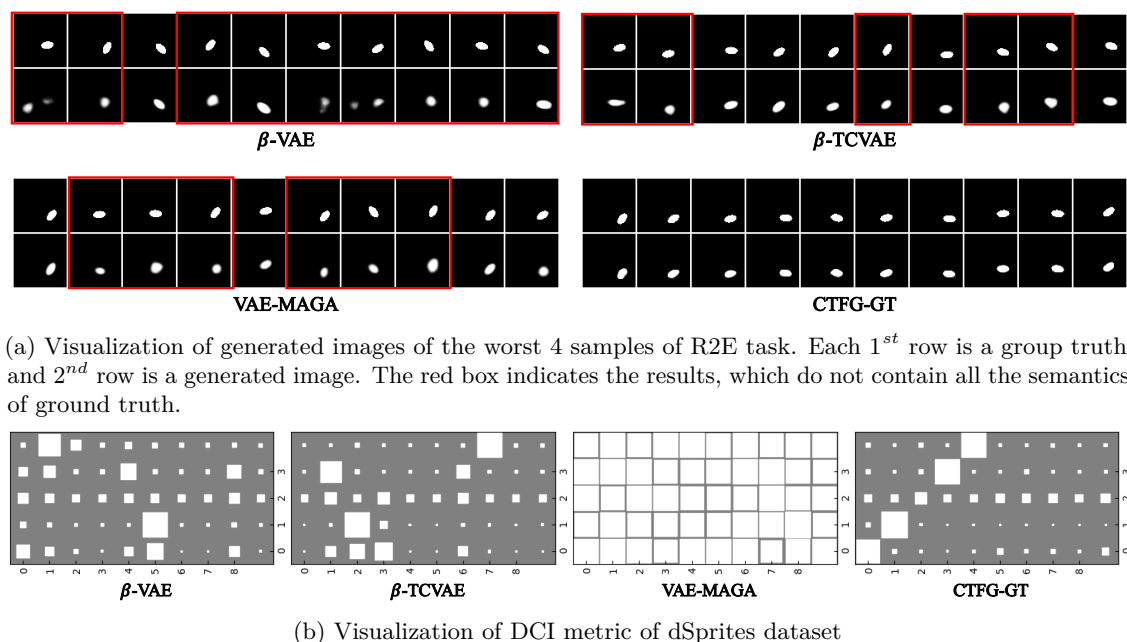
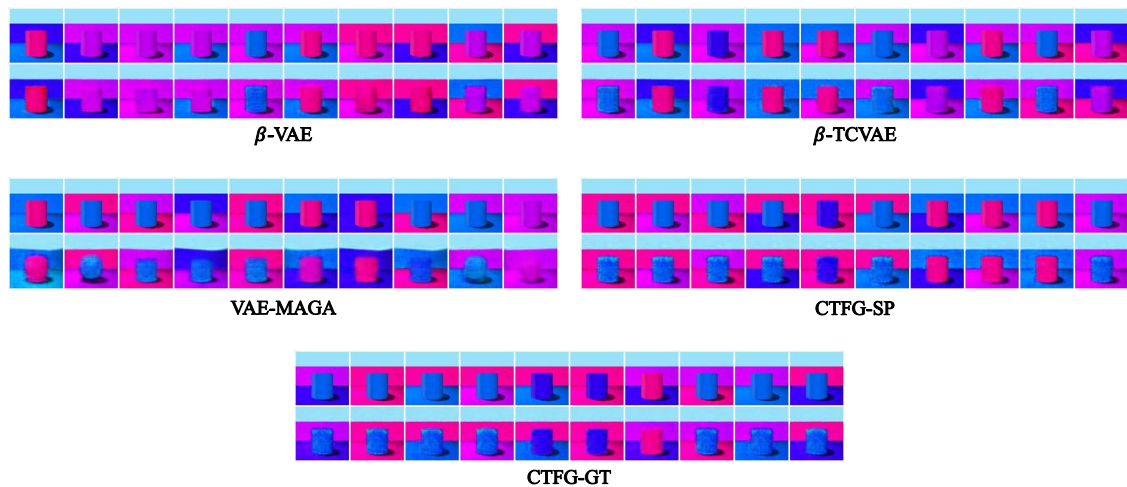
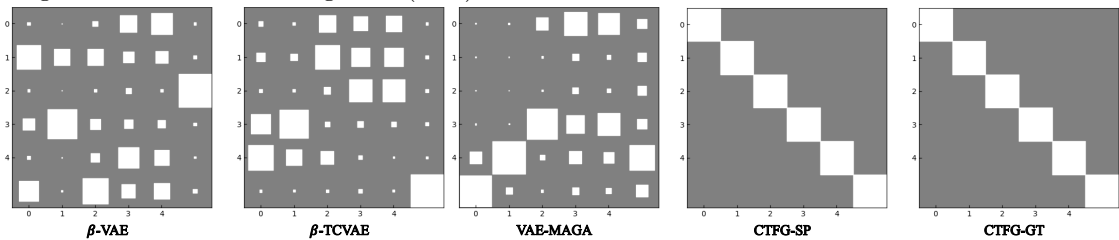


Figure 9: Qualitative results of compositional generalization of dSprites dataset (R2E). A more sparse matrix implies clear disentanglement.

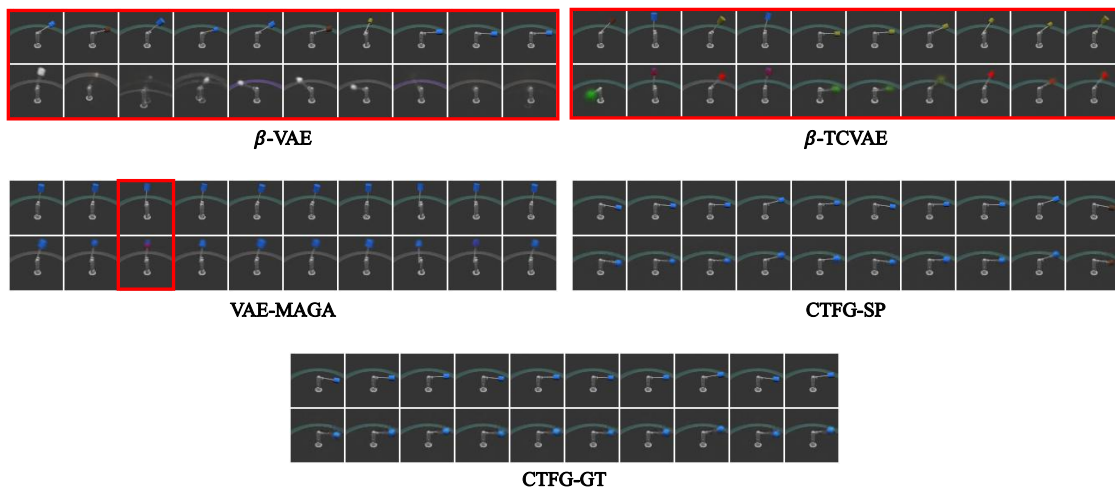


(a) Visualization of generated images of the worst 10 samples of the R2R task. Each 1<sup>st</sup> and 2<sup>nd</sup> row show the group truth samples and the generated results, respectively. The red box indicates the negative results, which do not contain all the semantics of the ground truth. We utilize randomly selected pivot images as introduced in Hwang et al. (2023) for the CTFG model.

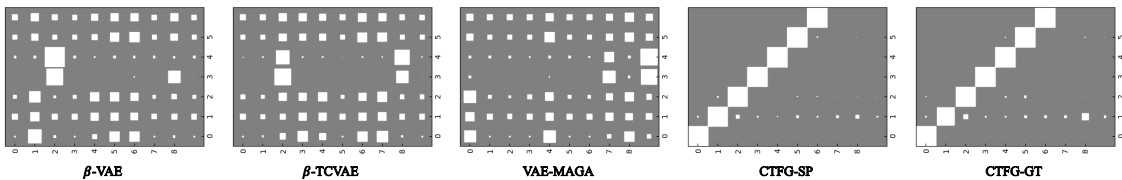


(b) Visualization of DCI metric of 3D Shapes dataset (training set). A more sparse matrix implies clear disentanglement.

Figure 10: Qualitative results of compositional generalization of 3D Shapes dataset (R2E).

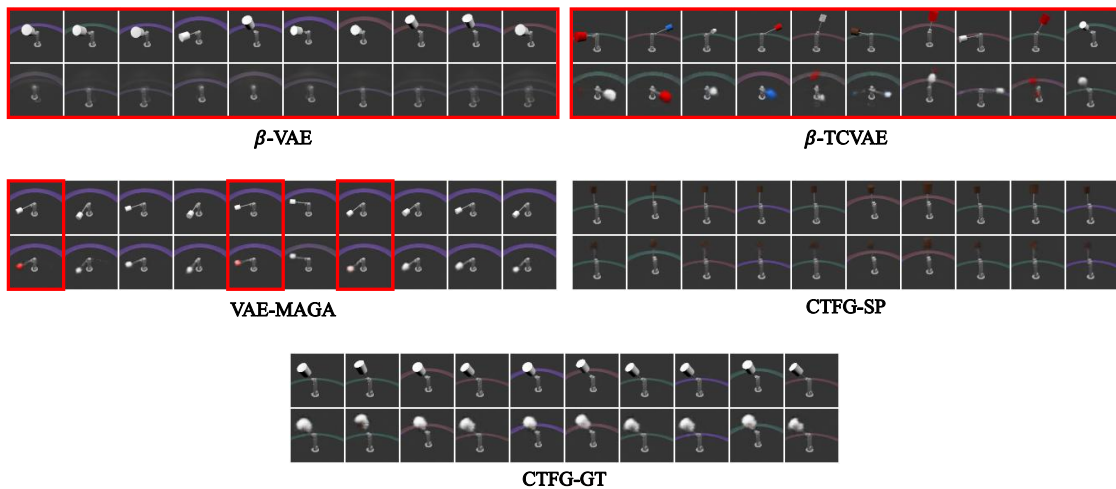


(a) Visualization of generated images of the worst 10 samples of the R2R task. Each 1<sup>st</sup> and 2<sup>nd</sup> row shows the group truth samples and the generated results, respectively. The red box indicates the negative results, which do not contain all the semantics of the ground truth. We utilize randomly selected pivot images as introduced in Hwang et al. (2023) for the CTFG model.

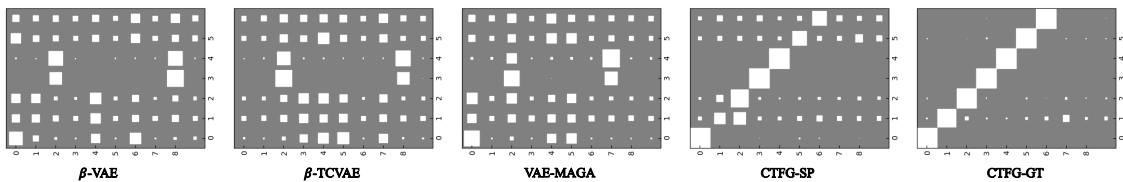


(b) Visualization of DCI metric of 3D Shapes dataset. A more sparse matrix implies clear disentanglement.

Figure 11: Qualitative results of compositional generalization of MPI3D dataset (R2E).



(a) Visualization of generated images of the worst 10 samples of the R2R task. Each 1<sup>st</sup> and 2<sup>nd</sup> row shows the group truth samples and the generated results, respectively. The red box indicates the negative results, which do not contain all the semantics of the ground truth. We utilize randomly selected pivot images as introduced in Hwang et al. (2023) for the CTFG model.



(b) Visualization of DCI metric of 3D Shapes dataset. A more sparse matrix implies clear disentanglement.

Figure 12: Qualitative results of compositional generalization of MPI3D dataset (R2R).

## F Disentanglement Learning Results

### F.1 Trade-Off (3D Shapes)

As illustrated in Figure 13, the proposed models improve the reconstruction error and disentanglement performance simultaneously on the dSprites dataset. Additionally, while the reconstruction error slightly decreases, the model performance dramatically improves compared to the baselines on the 3D Shapes dataset.

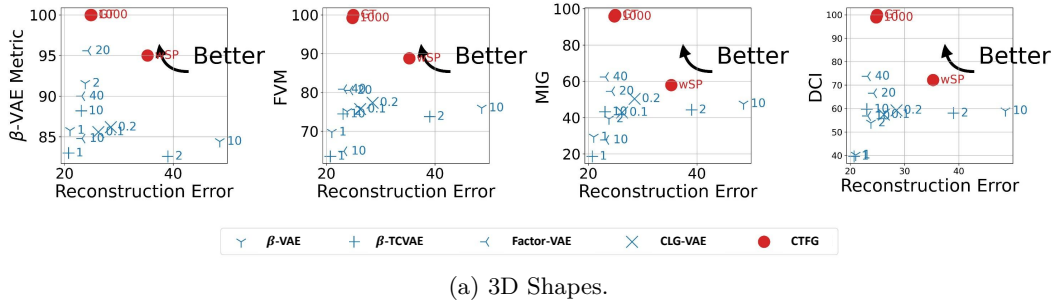


Figure 13: Reconstruction error vs. evaluation metrics ( $\beta$ -VAE metric, FVM, MIG, and DCI). The top left side indicates the best results for both objectives.

Table 6: Disentanglement performance of Homomorphism VAE vs. Ours with dSprites.

Model	beta-VAE	FVM	MIG	DCI
Homomorphism VAE	18.80( $\pm 5.75$ )	30.24( $\pm 12.18$ )	0.39( $\pm 0.76$ )	1.35( $\pm 1.12$ )
Groupified-VAE	79.30( $\pm 9.23$ )	69.75( $\pm 13.66$ )	21.03( $\pm 9.20$ )	31.08( $\pm 10.87$ )
CTFG-SP	<b>91.40</b> ( $\pm 4.99$ )	<b>93.74</b> ( $\pm 1.82$ )	<b>51.02</b> ( $\pm 2.42$ )	<b>64.69</b> ( $\pm 1.55$ )

### F.2 Ours vs. Homomorphism VAE

As shown in Table Table 6, our model outperforms Homomorphism VAE (Keurti et al., 2023) and Groupified-VAE (Yang et al., 2021). Homomorphism VAE (Keurti et al., 2023) utilizes the  $SO(2)$  for disentangled representation, and the elements of  $SO(2)$  affect the multi-dimension value of the latent vector. It implies that the  $SO(2)$  group is not appropriate symmetry for dimension-wise disentangled representation.

### F.3 Impact of Known Label Ratio

We set the known label ratio  $p \in \{0.1, 0.2, 0.4, 0.5\}$ . As shown in Table Table 7, our model is robust to the known label ratio except for the DCI metric.

Table 7: Disentanglement performance of weakly-supervised learning methods.

Model	dSprites				3DShapes			
	beta-VAE	FVM	MIG	DCI	beta-VAE	FVM	MIG	DCI
Ada-GVAE	83.60( $\pm 2.61$ )	83.67( $\pm 2.97$ )	21.34( $\pm 5.35$ )	<b>47.26</b> ( $\pm 1.89$ )	72.75( $\pm 6.50$ )	59.81( $\pm 6.14$ )	24.77( $\pm 7.48$ )	64.57( $\pm 4.04$ )
CTFG-wSP (0.1)	<b>88.60</b> ( $\pm 7.72$ )	83.36( $\pm 3.51$ )	14.71( $\pm 1.25$ )	23.46( $\pm 1.69$ )	<b>86.80</b> ( $\pm 3.90$ )	<b>84.05</b> ( $\pm 2.66$ )	<b>55.17</b> ( $\pm 2.18$ )	61.79( $\pm 3.15$ )
CTFG-wSP (0.2)	<b>89.78</b> ( $\pm 6.67$ )	<b>83.88</b> ( $\pm 2.76$ )	<b>23.03</b> ( $\pm 2.54$ )	28.07( $\pm 1.40$ )	<b>85.00</b> ( $\pm 13.11$ )	<b>83.00</b> ( $\pm 7.29$ )	<b>54.20</b> ( $\pm 13.35$ )	60.26( $\pm 12.00$ )
CTFG-wSP (0.4)	<b>88.20</b> ( $\pm 5.92$ )	83.49( $\pm 2.22$ )	<b>31.87</b> ( $\pm 2.02$ )	39.44( $\pm 0.79$ )	<b>86.40</b> ( $\pm 6.98$ )	<b>87.61</b> ( $\pm 7.09$ )	<b>61.47</b> ( $\pm 9.51$ )	<b>67.87</b> ( $\pm 8.39$ )
CTFG-wSP (0.5)	<b>87.00</b> ( $\pm 7.07$ )	<b>84.50</b> ( $\pm 1.41$ )	<b>31.95</b> ( $\pm 2.40$ )	39.36( $\pm 1.49$ )	<b>95.00</b> ( $\pm 7.07$ )	<b>88.81</b> ( $\pm 13.17$ )	<b>57.94</b> ( $\pm 16.52$ )	<b>72.14</b> ( $\pm 3.23$ )

### F.4 Impact of Each Loss

## G Qualitative Analysis of Disentanglement Learning

**3D Shapes** As shown in Figure 14, the baseline results show that multiple factors are changed when a single dimension value is changed. On the other hand, ours show the fully disentangled results represent:  $1^{st}$  row is the floor color changes,  $2^{nd}$  row is the wall color changes,  $3^{rd}$  row is the object color changes,  $4^{th}$  row is the scale of object,  $5^{th}$  row is the object shape changes, and  $6^{th}$  row is the rotation changes.

Table 8: Disentanglement performance over hyper-parameters

$(\alpha, \gamma)$	reconst. err.	beta-VAE	FVM	MIG	DCI
(100, 1)	<b>17.88</b> ( $\pm 1.24$ )	88.40( $\pm 4.97$ )	82.13( $\pm 1.86$ )	33.88( $\pm 3.03$ )	42.27( $\pm 2.02$ )
(200, 1)	21.23( $\pm 0.89$ )	92.44( $\pm 5.81$ )	87.54( $\pm 4.13$ )	35.21( $\pm 2.17$ )	47.40( $\pm 2.14$ )
(500, 1)	26.00( $\pm 2.72$ )	94.00( $\pm 4.42$ )	96.41( $\pm 1.83$ )	43.73( $\pm 3.55$ )	55.36( $\pm 2.70$ )
(1000, 1)	27.41( $\pm 0.73$ )	<b>95.80</b> ( $\pm 4.57$ )	<b>99.26</b> ( $\pm 1.12$ )	<b>51.81</b> ( $\pm 2.97$ )	<b>63.26</b> ( $\pm 2.73$ )

(a) dSprites with Ground Truth model

$(\alpha, \beta, \gamma)$	reconst. err.	beta-VAE	FVM	MIG	DCI
(100, 1, 1)	<b>19.53</b> ( $\pm 1.75$ )	85.78( $\pm 5.87$ )	82.10( $\pm 2.81$ )	27.12( $\pm 1.98$ )	34.42( $\pm 1.04$ )
(100, 2, 1)	17.58( $\pm 1.08$ )	86.44( $\pm 5.90$ )	82.40( $\pm 1.62$ )	34.08( $\pm 1.75$ )	41.16( $\pm 0.98$ )
(1000, 1, 1)	26.22( $\pm 1.31$ )	<b>91.40</b> ( $\pm 4.90$ )	93.46( $\pm 1.97$ )	46.35( $\pm 1.94$ )	57.92( $\pm 1.97$ )
(1000, 2, 1)	26.32( $\pm 2.20$ )	<b>91.40</b> ( $\pm 4.99$ )	<b>93.74</b> ( $\pm 1.82$ )	<b>51.03</b> ( $\pm 2.42$ )	<b>64.69</b> ( $\pm 1.55$ )

(b) dSprites with Supervised method

$(\alpha, \gamma)$	reconst. err.	beta-VAE	FVM	MIG	DCI
(100, 1)	33.62( $\pm 5.38$ )	89.60( $\pm 6.17$ )	82.44( $\pm 3.78$ )	53.37( $\pm 12.87$ )	60.04( $\pm 13.98$ )
(200, 1)	34.50( $\pm 5.35$ )	95.00( $\pm 5.34$ )	91.95( $\pm 7.14$ )	68.06( $\pm 17.82$ )	74.16( $\pm 14.69$ )
(500, 1)	29.30( $\pm 1.72$ )	<b>100.00</b> ( $\pm 0.00$ )	97.28( $\pm 3.05$ )	86.68( $\pm 5.56$ )	90.68( $\pm 5.87$ )
(1000, 1)	<b>24.94</b> ( $\pm 1.51$ )	<b>100.00</b> ( $\pm 0.00$ )	<b>100.00</b> ( $\pm 0.00$ )	<b>95.57</b> ( $\pm 0.80$ )	<b>99.94</b> ( $\pm 0.18$ )

(c) 3D Shapes with Ground Truth model

**MPI3D** As shown in Figure 15, the baseline results show that multiple factors are changed when a single dimension value is changed. Also, the object usually disappeared following the intervals with baselines. On the other hand, supervised methods show better results than baselines: 1<sup>st</sup> row is the object color changes, 2<sup>nd</sup> row is the object shape changes, 3<sup>rd</sup> row is the object size changes, 5<sup>th</sup> row is the background color changes, 7<sup>th</sup> row is the vertical axis changes, and 9<sup>th</sup> row is the height changes. Also, the GT model results represent: 1<sup>st</sup> row is the object color changes, 2<sup>nd</sup> row is the object shape changes, 3<sup>rd</sup> row is the object size changes, 4<sup>th</sup> row is the height changes, 5<sup>th</sup> row is the background color changes, 6<sup>th</sup> row is the vertical axis changes, and 7<sup>th</sup> row is the horizontal axis changes.

**cLPR** As shown in Figure 16, the Homomorphism VAE (Keurti et al., 2023) shows that multiple factors are changed when a single dimension value is changed. Also, the reconstruction quality is lower than ours (CTFG-GT and CTFG-SP). On the other hand, the supervised method shows better results than the Homomorphism VAE: 1<sup>st</sup> and 3<sup>rd</sup> rows are z-axis rotation, 2<sup>nd</sup> row is y-axis rotation, and 4<sup>th</sup> row is x-axis rotation. Also, the GT model results represent: 1<sup>st</sup> and 2<sup>nd</sup> rows are x-axis rotation, 3<sup>rd</sup> row is z-axis rotation, and 4<sup>th</sup> and 6<sup>th</sup> rows are y-axis rotation.

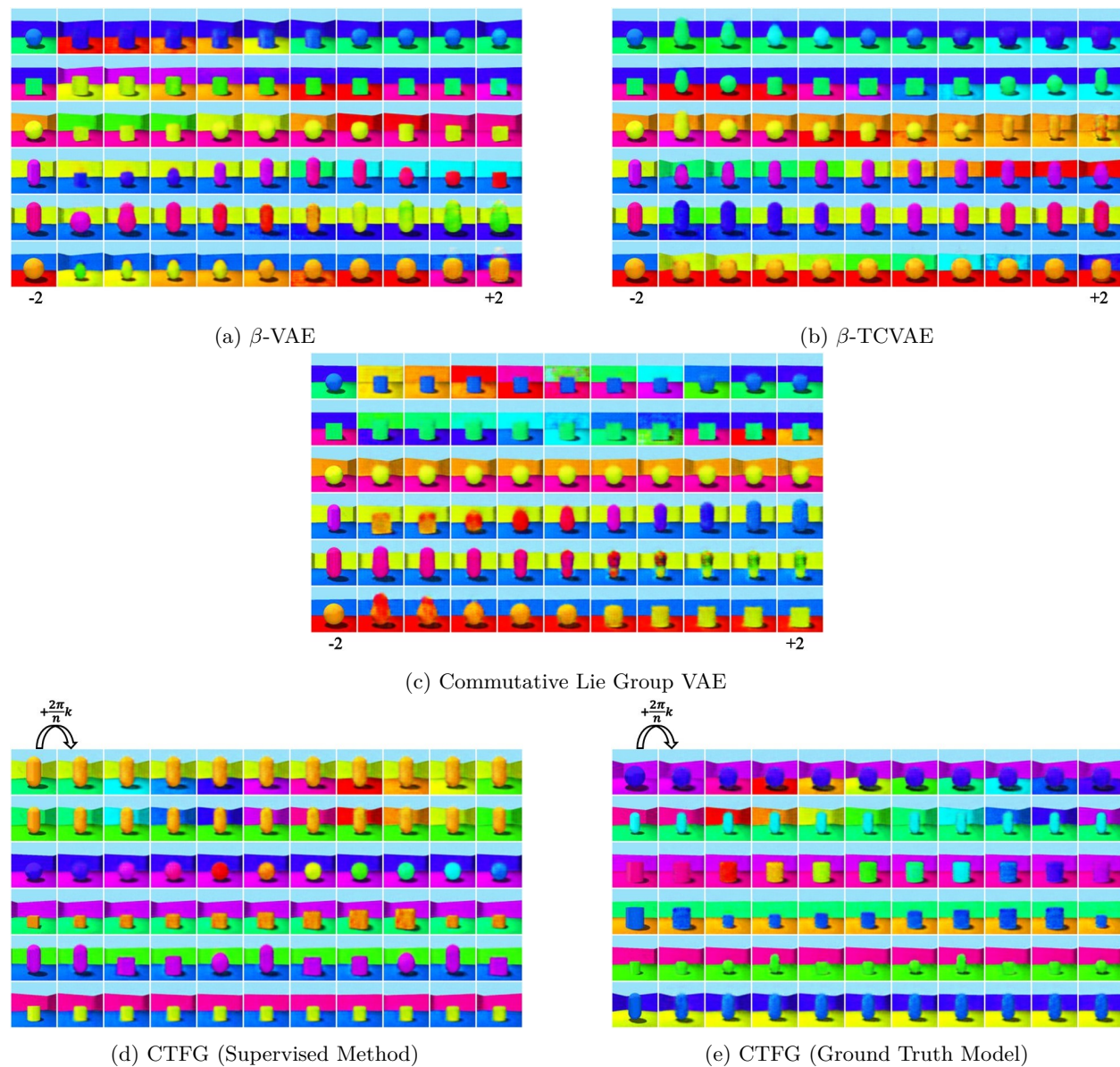


Figure 14: The 1<sup>st</sup> column images are randomly selected from the dataset. Each row indicates each dimension of each model.  $\beta$ -VAE,  $\beta$ -TCVAE, and Commutative Lie Group VAE trace each dimension value from -2 to +2. The proposed methods apply a group action  $+\frac{2\pi}{n}$  to the selected images a total of 10 times.

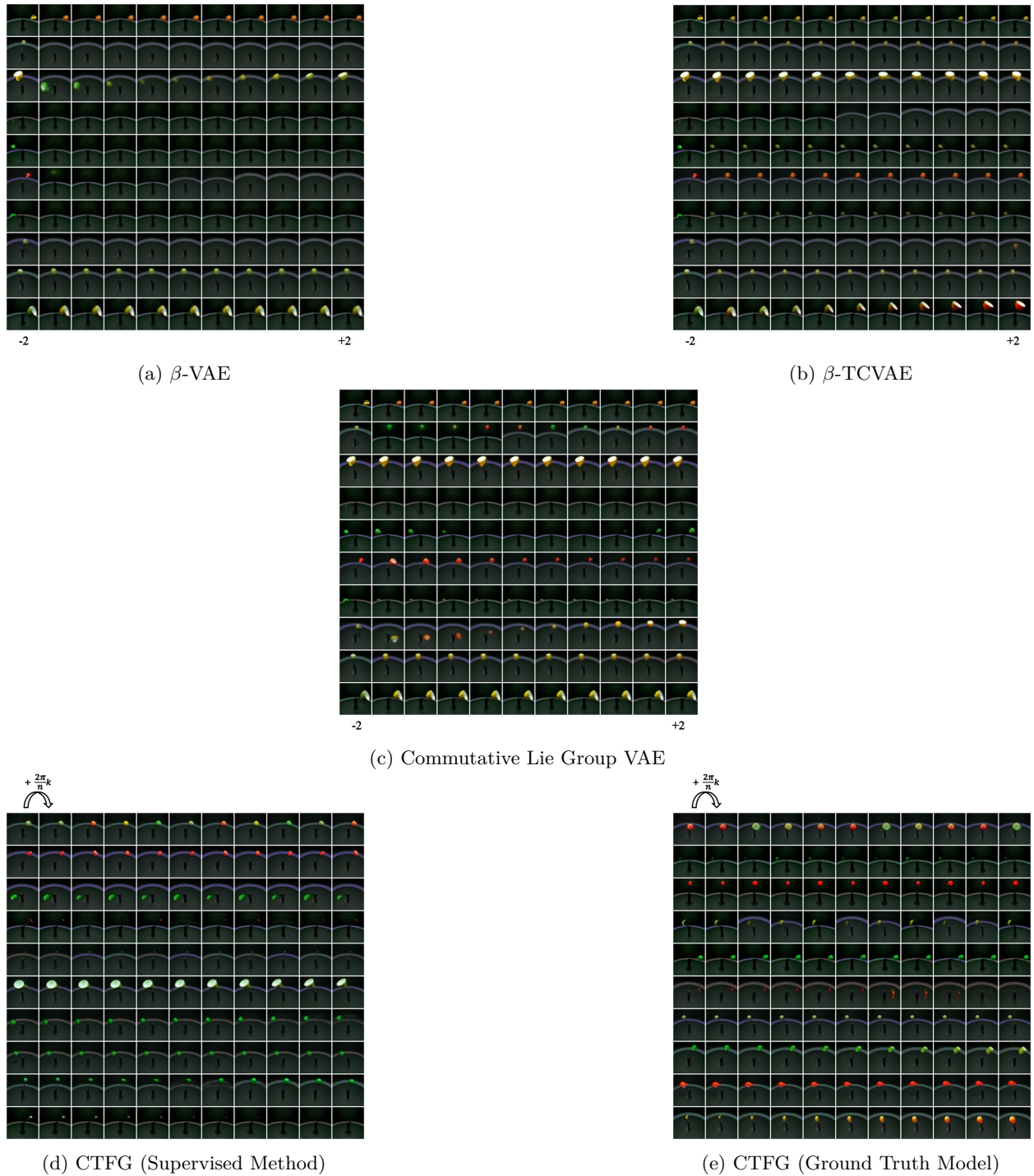
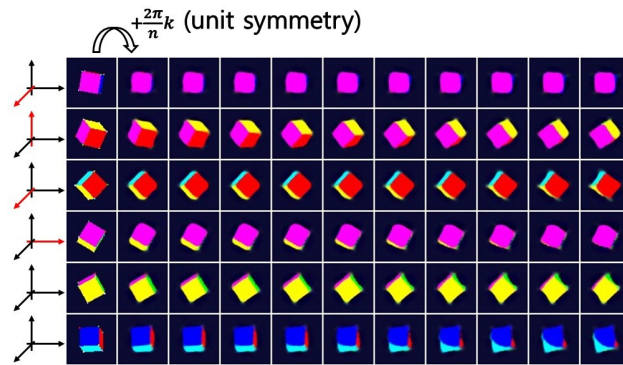
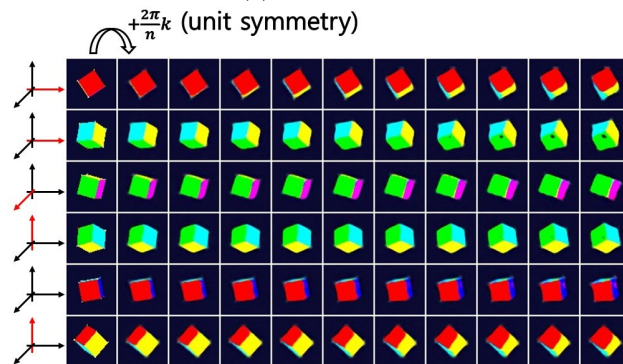


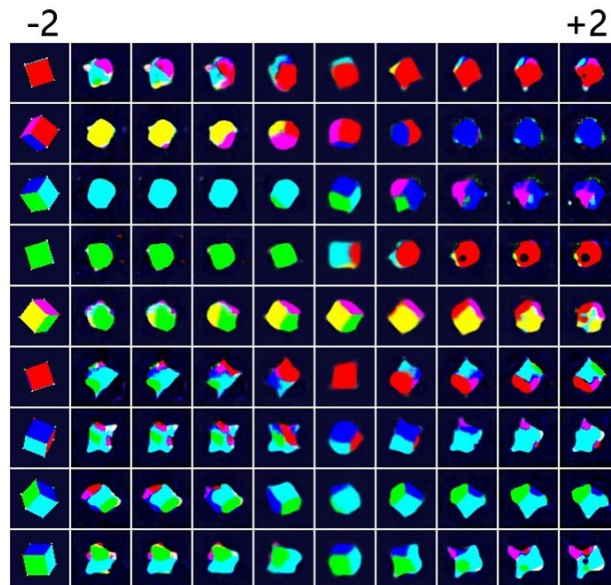
Figure 15: The 1<sup>st</sup> column images are randomly selected from the dataset. Each row indicates each dimension of each model.  $\beta$ -VAE,  $\beta$ -TCVAE, and Commutative Lie Group VAE trace each dimension value from -2 to +2. The proposed methods apply a group action  $+\frac{2\pi}{n}$  to the selected images a total of 10 times.



(a) CTFG-GT



(b) CTFG-Super



(c) Homomorphism VAE

Figure 16: The 1<sup>st</sup> column images are randomly selected from the dataset. Each row indicates each dimension of each model. CTFG-GT, CTFG-Super ( $\alpha: 100.0$ ), and homomorphism VAE trace each dimension value from -2 to +2. The proposed methods apply a group action  $+\frac{2\pi}{n}$  to the selected images a total of 10 times. And red color axis is a rotation axis.