

SELF-SUPERVISED DISENTANGLEMENT VIA CLUSTER-DEPENDENT ROTATIONAL EQUIVARIANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Conventional self-supervised learning methods extract robust features by enforcing invariance to data augmentations. While effective for obtaining clustered representations, this objective provides limited control over how data variations structure the feature space, hindering disentanglement. Recent methods improve feature space structure by imposing equivariant predictability on feature transformations induced by data augmentations. However, existing approaches suffer from two significant limitations: (i) the incorporation of invariance in their final objective interferes with the learning of neat equivariance; (ii) the imposition of uniform equivariance across all samples forces semantic clusters into a parallel arrangement, leading to reduced inter-cluster distances (for features on the hypersphere). To overcome these issues, we propose in this paper Cluster-Dependent Rotational Equivariance for Disentanglement (CD-RED), a framework that enables learning neat equivariance and uniformly distributed clusters, while further supporting perfect disentanglement. Notably, CD-RED explicitly encodes variations as rotations via a direct product of $SO(2)$ groups within orthogonal hyperspherical subspaces, providing a principled mechanism for precise equivariance. We theoretically and experimentally establish that CD-RED achieves perfectly disentangled representations, suggesting a promising new direction for self-supervised disentanglement.

1 INTRODUCTION

Learning disentangled representations, where independent factors of data variation are encoded into separate feature dimensions, is a fundamental goal in machine learning (Bengio et al., 2013; Hinton & Salakhutdinov, 2006; Higgins et al., 2018; Locatello et al., 2019; Weiler & Cesa, 2019). Such representations are critical for a wide range of purposes, including improving generalization, enhancing interpretability, data generation, and transfer learning (Locatello et al., 2019; Ren et al., 2021).

The pursuit of disentanglement has historically followed two main paradigms: supervised learning with labeled factors (Reed et al., 2014; Cheung et al., 2014), and unsupervised learning with tailored inductive biases designed to encourage statistical independence in the latent space (Higgins et al., 2017; Kim & Mnih, 2018; Chen et al., 2018). The former is limited by the cost of labels, while the latter often relies on strong assumptions that may not align with the real-world data (Locatello et al., 2019). This landscape has motivated a recent shift towards self-supervised learning, which seeks to discover factors directly from the structure of unlabeled data.

Conventional self-supervised learning methods (Tian et al., 2020; He et al., 2020; Wu et al., 2018; Chen et al., 2020a) extract robust features by enforcing invariance to data augmentations, typically by minimizing a contrastive loss such as InfoNCE (Oord et al., 2018). Despite their popularity and effectiveness for clustering, these methods often fail to achieve feature disentanglement in practice, since they lack precise control over how the learned representation structures features associated with data variations — merely requiring proximity. Although some analyses suggest InfoNCE can promote disentanglement under strong assumptions (Von Kügelgen et al., 2021; Zimmermann et al., 2021; Ngweta et al., 2023), the conditions rely on idealizations that are seldom met in practice.

To impose more structure, one line of work (Xiao et al., 2020; Wang et al., 2021; Eastwood et al., 2023) stacks a group of contrastive objectives, working on different partitions of the data or working with distinct sets of augmentations, to extract factors of data variation. While effective to some

054 degree, these methods exhibit high conceptual complexity and still typically rely on strong assump-
 055 tions or independence regularizations to achieve disentanglement.

056 Another prominent line of work seeks to go beyond mere invariance by learning equivariant rep-
 057 resentations (Dangovski et al., 2021), where feature transformations predictably correspond to data
 058 transformations. While this provides a principled way to organize features towards disentanglement,
 059 existing equivariant methods face significant limitations. For instance, the approach by Shakerinava
 060 et al. (2022) is limited to single-cluster manifolds. Garrido et al. (2023); Marchetti et al. (2023) rely
 061 on modeling cluster and variation features in separate spaces.

062 Among these distinct approaches, our work follows Gupta et al. (2023): learning representations that
 063 capture discrete latent structures on a single hypersphere. Nevertheless, they require jointly optimiz-
 064 ing the InfoNCE and equivariance losses. This coupling creates an inherent trade-off that prevents
 065 neat equivariance and, consequently, equivariance-based disentanglement. To address this issue, we
 066 propose a two-stage framework in this paper that confines the use of InfoNCE to the initialization
 067 stage, which allows the second stage to focus solely on equivariance learning. We thereby achieve
 068 neat equivariance, a fundamental prerequisite for equivariance-based disentanglement.

069 However, neat equivariance alone is insufficient for disentanglement, as the geometric structure of
 070 the equivariance is also critical Shakerinava et al. (2022). A common issue in existing methods is
 071 that they employ an overly arbitrary equivariance modeling, which often fails to impose a meaning-
 072 ful geometric transformation. For instance, Devillers & Lefort (2022) uses a learnable feed-forward
 073 network to predict feature changes, which cannot guarantee a consistent geometric structure. In con-
 074 trast, Gupta et al. (2023) adopts a more structured approach by modeling the changes as an implicit
 075 rotation on the hypersphere. Although this approach has clear potential to structure the feature space,
 076 a significant limitation is its mandate of a global transformation that is identical for all inputs. This
 077 inherently forces the semantic clusters into a parallel arrangement, thereby reducing inter-cluster
 078 distances and, consequently, the representational efficiency on the hypersphere. To overcome this
 079 limitation, we propose in this paper Cluster-Dependent Rotational Equivariance, which resolves the
 080 issue by making the rotational transformation local to each cluster.

081 Another key challenge in learning neat equivariance lies in the imprecise nature of implicit modeling.
 082 To address this, we shift to an explicit model of rotational transformation, leveraging a composite
 083 rotation group formed by the direct product of $SO(2)$ subgroups (Quessard et al., 2020; Shakerinava
 084 et al., 2022) to cleanly encode data variations as rotations within a set of orthogonal 2D subspaces.

085 Integrating these insights, we arrive at Cluster-Dependent Rotational Equivariance for Disentangle-
 086 ment (CD-RED). Our theoretical and experimental analyses show that CD-RED achieves perfectly
 087 disentangled representations, constituting a significant advance in self-supervised learning.

088 Our primary contributions are summarized as follows:

- 089 • A novel two-stage learning framework that decouples invariance and equivariance learning, en-
 090 abling learning neat equivariant representations on a single hypersphere.
- 091 • The introduction of Cluster-Dependent Rotational Equivariance that overcomes the inherent lim-
 092 itations of global rotational equivariance by making the transformation local to each cluster.
- 093 • CD-RED, a concrete method that provably achieves perfect disentangled representations.
- 094 • A theoretical extension of disentanglement theory, moving beyond the idealized assumptions of
 095 Higgins’ formulation to handle more practical, complex augmentations.

099 2 PRELIMINARY AND RELATED WORK

100 We begin by formalizing the representation learning pipeline and clarifying how invariance, equivari-
 101 ance, and disentanglement arise within this context.

102 We adopt the standard generative formulation where data originates from a latent semantic space
 103 $\mathcal{Z} \in \mathbb{R}^m$ composed of independent semantic factors (z_1, \dots, z_m) . Each $z \in \mathcal{Z}$ is rendered into the
 104 observation space \mathcal{X} via a possibly unknown invertible and non-linear function $\varphi : \mathcal{Z} \mapsto \mathcal{X}$, and
 105 mapped into a feature space by a **feature extractor** $f : \mathcal{X} \mapsto \mathcal{Y}$, giving rise to the representation
 106 pipeline: $\Phi : \mathcal{Z} \xrightarrow{\varphi} \mathcal{X} \xrightarrow{f} \mathcal{Y}$. We assume $f(x)$ is ℓ_2 -normalized unless otherwise stated.
 107

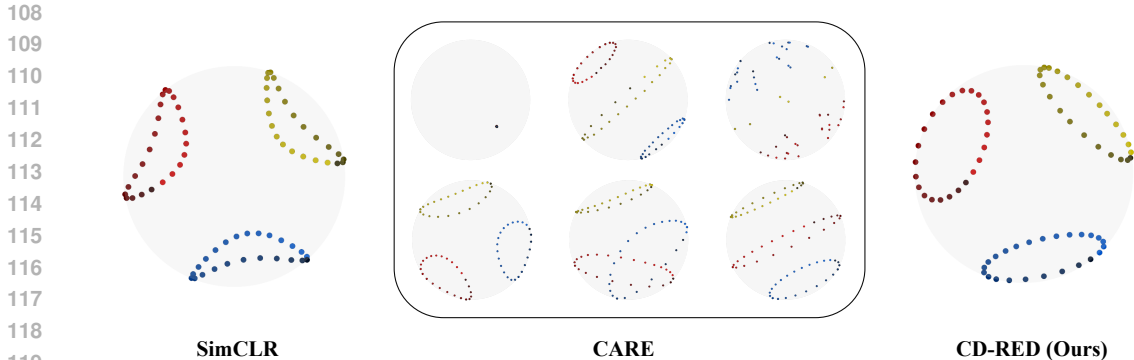


Figure 1: An Illustration of the essential differences between closely related methods with results on a dataset of three rotating objects. **SimCLR** (Chen et al. (2020a), InfoNCE): without equivariance, the learned features are not well-structured, though forming uniformly distributed clusters. **CARE** (Gupta et al., 2023): jointly optimizing the InfoNCE and global rotational equivariance losses, the results are highly sensitive to the weighting of each component; under a good balance, the learned features form parallelly distributed clusters with varying radius. **CD-RED**: with cluster-dependent rotational equivariance, the learned features form uniformly distributed, well-structured clusters.

To enable training in the self-supervised setting, a set of m data augmentations $\mathcal{A} = \{a_i\}_{i=1}^m$ is considered, where each a_i denotes a transformation (e.g., rotation, cropping), possibly parameterized by a latent variable t drawn from T_{a_i} . Each a_i can be applied individually or in combination with others. For simplicity, we abuse $a \in \mathcal{A}$ as a single or composite augmentation, while retaining the ability to refer to internal parameters t when needed.

Given an input $x = \varphi(z)$, we construct a positive sample $x^+ = a(x)$ by applying a sampled augmentation $a \in \mathcal{A}$, and aim to relate $f(x)$ and $f(x^+)$ through appropriate learning objectives.

A core setting is that each augmentation in \mathcal{A} preserves some latent factors while perturbing others. We categorize the latent space \mathcal{Z} into two subsets¹:

- **Style Latents**: Latent factors that are altered by at least one augmentation in \mathcal{A} , which encodes changes such as views, colors, and positions.
- **Content Latents**: Latent factors that are preserved under all augmentations in \mathcal{A} . These define the semantic identity of the sample and are expected to remain unchanged under transformation.

This distinction underpins the goals of self-supervised learning: to capture content-related variations through invariance, while optionally modeling style-related variations through equivariance.

Augmentations in Invariance A foundational goal in self-supervised learning is to learn invariant representations, i.e., feature embeddings that remain stable under semantic-preserving augmentations. The objective can be characterized by $f(a(x)) \approx f(x), \forall a \in \mathcal{A}$, while an alternative formula can be $f(a_1(x)) \approx f(a_2(x)), \forall a_1, a_2 \in \mathcal{A}$, which is primarily the same. We would express it in terms of the former for simplicity. A widely-used objective for learning such invariant representations is the InfoNCE (Oord et al., 2018; Tian et al., 2020; He et al., 2020; Wu et al., 2018; Chen et al., 2020a), which is formulated as:

$$\mathcal{L}_{\text{Info}} = \frac{1}{2N} \sum_{i=1}^{2N} -\log \frac{\exp(\text{sim}(f(x_i), f(x_i^+))/\tau)}{\sum_{j=1, j \neq i}^{2N} \exp(\text{sim}(f(x_i), f(x_j))/\tau)}, \quad (1)$$

where the features $f(x)$ are confined to the hypersphere $\mathcal{Y} = \mathbb{S}^{d-1} \triangleq \{z \in \mathbb{R}^d : \|z\|_2 = 1\}$, $\text{sim}(a, b) = a^\top b$ denotes the cosine similarity between two unit vectors, and $\tau > 0$ is the temperature parameter. This loss encourages features of augmented views of the same semantic instance to cluster, while driving features of different instances to be uniformly distributed on the hypersphere (Parulekar et al., 2023; Wang et al., 2022; Wang & Isola, 2020).

¹Note that the assignment of latents as content or style is contingent upon the chosen augmentations.

Augmentations in Equivariance In contrast to invariance, equivariance structures the feature space by requiring that transformations in the input space induce predictable changes in the feature space, which can be formulated as $f(a(x)) \approx M_a(f(x))$, where M_a is a feature-space transformation aligned with a . The corresponding objective can be:

$$\mathcal{L}_{\text{equi}_0} = \mathbb{E}_{x,a} [\|f(a(x)) - M_a(f(x))\|^2]. \quad (2)$$

Optimizing this objective alone may lead to collapsed features, and this objective is therefore usually optimized jointly with a collapse-avoiding loss (e.g., InfoNCE). Existing methods differ in their modeling of M_a . One class of methods learns M_a as a neural network conditioned on both $f(x)$ and a , i.e., $M_a(f(x)) = \text{MLP}(f(x), a)$ (Xiao et al., 2020; Dangovski et al., 2021; Devillers & Lefort, 2022). Predicting change with an arbitrarily expressive non-linear neural network can enable the model to learn richer and more informative features; however, it may contribute little to structuring the feature space. In contrast, more recent methods impose stronger geometric structures. For instance, Garrido et al. (2023) models the transformation more structurally as parametrized linear mappings, while (Shakerinava et al., 2022; Gupta et al., 2023) implicitly constrain M_a to be, e.g., an orthogonal matrix, giving the loss as $\mathcal{L}_{\text{equi}_1} = \mathbb{E}_{a \in \mathcal{A}} \mathbb{E}_{x, x' \in \mathcal{X}} [f(a(x'))^\top f(a(x)) - f(x)^\top f(x')]^2$. In fact, for features constrained to a hypersphere, these two approaches are nearly equivalent, since any valid linear transformation operating on a hypersphere reduces to an orthogonal matrix. Notably, they enforce such orthogonal transformation globally across all samples on the hypersphere, which drives features of distinct semantic clusters to be parallelly distributed on different hyperplanes, conflicting with the commonly coupled InfoNCE objective that promotes uniformly distributed clusters across the hypersphere and, more critically, inherently leading to reducing inter-cluster distances.

Augmentations in disentanglement The objective of disentangled representation learning is to induce a feature mapping f such that independent latent factors of variation are encoded into distinct, interpretable components within the learned representation $f(x)$ (Bengio et al., 2013). A foundational framework was put forward by Higgins et al. (2018), which defines a representation as disentangled if variations in a single generative factor correspond to changes in a single latent unit, with all other units remaining unaffected. Building on this conceptual foundation, recent self-supervised methods have aimed to approximate disentanglement through group-theoretic formulations (Wang et al., 2021; Marchetti et al., 2023; Shakerinava et al., 2022). In these approaches, data augmentations are assumed to form a decomposable group $\mathcal{G} = \mathcal{G}_1 \times \dots \times \mathcal{G}_n$, with each subgroup \mathcal{G}_i acting on a single generative factor and corresponding to a specific subspace \mathcal{Y}_i in the feature space. While elegant in theory, this framework hinges on strong assumptions about the augmentations that each affects a single generative factor. In practice, however, many augmentations simultaneously affect multiple latent factors. In contrast to prior work, we also model augmentations that simultaneously impact multiple latent factors, enabling disentanglement under such entangled augmentations.

3 CLUSTER-DEPENDENT ROTATIONAL EQUIVARIANCE

We introduce in this section the proposed framework, Cluster-Dependent Rotational Equivariant for self-supervised Disentangled representation learning (CD-RED).

Our first insight regarding cluster-dependent rotational equivariance is that while equivariance is essential, its joint optimization with a collapse-preventing objective like InfoNCE inherently interferes with learning neat equivariance. However, the core purpose of the collapse-preventing objective is solely to avoid feature collapse, which can be better achieved with a two-stage method to avoid interference: a primary stage using self-supervised clustering (e.g., InfoNCE) to establish non-collapsed semantic clustering, followed by a secondary stage that learns neat equivariance within each cluster, where ensuring non-trivial equivariance suffices to prevent intra-cluster collapse.

Our second insight regarding cluster-dependent rotational equivariance is that enforcing global rotational equivariance (e.g., that employed by CARE, Gupta et al. (2023)) for features on the hypersphere would drive distinct semantic clusters to be parallelly distributed on different hyperplanes. This inherently results in reduced inter-cluster distances, blurs the boundaries between clusters, and ultimately limits the framework to supporting the joint modeling of only a relatively small number of clusters. This limitation can be naturally addressed by shifting to a cluster-dependent rotational system, i.e., making the transformation local to each cluster.

We summarize our method in Algorithm 1, outlining the overall objective and the training steps.

3.1 STAGE 1: SELF-SUPERVISED INVARIANCE LEARNING FOR CLUSTERING

In the first stage, we perform conventional self-supervised invariance learning by the InfoNCE loss to obtain well-structured representations. Theoretically, at its optimum, InfoNCE learns semantically meaningful clusters that are uniformly distributed across the hypersphere (Wang & Isola, 2020).

This initial stage yields features that are amenable to clustering. We then apply spherical k-means (Hornik et al., 2012) or agglomerative clustering (Müllner, 2011) to these features to obtain the initial cluster centroids and assignments.

3.2 STAGE 2: SELF-SUPERVISED EQUIVARIANCE LEARNING FOR DISENTANGLEMENT

In the second stage, we learn cluster-dependent rotational equivariance within each cluster.

3.2.1 CLUSTER-DEPENDENT ROTATION SYSTEMS

Cluster-dependent rotational equivariance necessitates a local rotation system for each cluster. This requires two key components: a rotation center and a rotation group.

Rotation Center. For the rotation center, we initialize it as the cluster centroid from the first stage, which serves as the anchor point for all rotations within the cluster. Concretely, each cluster i is assigned a rotation center r_i , defined as the normalized centroid of its features: $r_i = c_i$. Since InfoNCE is proven to induce uniformly distributed clusters, these initial centroids already provide stable rotation anchors and can thus be fixed during subsequent training. Nevertheless, we may optionally introduce the loss in Eq. (31) to further promote its uniform distribution.

Rotation Group. Existing approaches often model equivariant transformations as generalized linear transformations. When applied to features constrained to hypersphere surfaces, they reduce to orthogonal matrices (Gupta et al., 2023). This, however, introduces two main issues.

- *Reflection ambiguity:* Orthogonal matrices may include reflections ($\det = -1$), which can flip semantic orientations. While the model may not learn reflections if they do not fit the data, this remains a potential drawback.
- *Unnecessary complexity:* Full $d \times d$ rotation involves many feature dimensions, making it unnecessarily complex for modeling simple augmentations that typically affect low-dimensional subspaces. Such complexity not only increases the learning difficulty but also hampers subsequent disentanglement.

To address these issues, we restrict each transformation to a 2D subspace and model it directly via an element of the special orthogonal group $SO(2)$, defined as

$$SO(2) \triangleq \{R \in \mathbb{R}^{2 \times 2} \mid R^\top R = I, \det(R) = 1\}. \quad (3)$$

Each element of $SO(2)$ can be parameterized by a single angle $\theta \in [-\pi, \pi]$ via the mapping

$$R(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}. \quad (4)$$

We assign each orthogonal augmentation² to a dedicated 2D hyperspherical subspace within the full feature space \mathbb{S}^{d-1} and restrict the corresponding feature rotations to that subspace. More formally, for the i -th augmentation, we assign it to operate on the $(2i-1)$ -th and $(2i)$ -th dimensions without loss of generality: each orthogonal augmentation must occupy at least two dimensions, but it does not matter which ones. Instead of letting the model learn to select these dimensions, we specify them directly, simplifying the process and avoiding unnecessary complexity. Specifically, for n orthogonal augmentations, we define the composite rotation as a block-diagonal matrix:

$$\mathcal{R}(\theta) = \text{diag}(R(\theta^{(1)}), R(\theta^{(2)}), \dots, R(\theta^{(n)}), I_{d-2n}) \in \mathbb{R}^{d \times d}, \quad (5)$$

which subsequently acts on the full feature vector.

Each θ_i is predicted by a small **neural network** g_i dedicated to the i -th augmentation. The input to g_i can include the augmentation parameter t (e.g., rotation angle) and, optionally, input-dependent

²Augmentations that cause variations in a single style factor that are distinct from each other.

signals (for input-dependent augmentations). The output is constrained to $[-\pi, \pi]$ via a tanh activation to: (i) avoid degenerate 2π -periodic equivariance, and (ii) enable inverse operations (e.g., left vs. right shift) by supporting both positive and negative values.

This construction not only avoids reflections and unnecessary complexity but also provides a clear, modular way to encode independent variation directions. It naturally generalizes to more complex cases: for augmentation that may affect multiple latent factors, we assign it to the multiple orthogonal subspaces and learn corresponding rotation parameters (see Section 4.3).

Coordinate Alignment. The block-diagonal rotation matrix $\mathcal{R}(\theta)$ introduced earlier operates in a fixed coordinate system. To accommodate cluster-dependent centers, we align the rotation center of each cluster to a canonical reference direction, which is chosen as the north pole $e_d = (0, 0, \dots, 1) \in \mathbb{S}^{d-1}$, before applying the rotation. This ensures that all rotations act on aligned coordinates while preserving cluster-local dynamics.

We achieve this alignment using the *Householder transformation* (Householder, 1958), a reflection over a hyperplane orthogonal to a given vector u :

$$H = I - 2 \frac{uu^\top}{u^\top u}, \quad (6)$$

where u is the normal vector of the reflecting hyperplane. In particular, a vector $v \in \mathbb{S}^{d-1}$ can be reflected to another vector $w \in \mathbb{S}^{d-1}$, i.e., $Hv = w$, by setting let $u = v - w$. To reflect the rotation center r_i of cluster i to e_d , we set $u_i = r_i - e_d$.

3.2.2 ENFORCING EQUIVARIANCE WITH ALIGNED ROTATIONS

With the constructed rotation systems, our equivariance condition can be expressed as:

$$H_{c(x)}f(a(t, x)) = \mathcal{R}(\theta_{a(t)})H_{c(x)}f(x), \quad \forall x \in \mathcal{X}, a(t) \in \mathcal{A}, \quad (7)$$

where $c(x) := \arg \max_i f(x)^\top r_i$ denotes the dynamically reassigned cluster of x determined by its proximity to the rotation centers, $H_{c(x)}$ is the corresponding Householder matrix, and $\mathcal{R}(\theta_{a(t)})$ is the block-rotation matrix for augmentation $a(t)$, defined as

$$\mathcal{R}(\theta_{a(t)}) := \text{diag}(R(\theta_{a(t)}^{(1)}), \dots, R(\theta_{a(t)}^{(m)}), I_{d-2m}), \quad (8)$$

with $\theta_{a(t)}^{(i)}$ being the rotation angle predicted by the corresponding g network for augmentation a in the i -th subspace. For composite augmentation, $\theta_{a(t)}^{(i)}$ is the accumulated rotation angle over the constituent single augmentations. We decompose the coordinate-aligned feature $H_{c(x)}f(x)$ for subsequent use:

$$H_{c(x)}f(x) = ([H_{c(x)}f(x)]_1, \dots, [H_{c(x)}f(x)]_m, [H_{c(x)}f(x)]_{\text{extra}}), \quad (9)$$

where each subspace block $[\cdot]_i \in \mathbb{R}^2$ and $[\cdot]_{\text{extra}} \in \mathbb{R}^{d-2m}$.

Equivariance is prone to trivial solutions. For instance, $\mathcal{R}(\theta_{a(t)})$ may collapse to the identity, yielding $f(a(x)) = f(x)$ and thus reducing to invariance rather than true equivariance. Existing remedies, such as combining with an InfoNCE loss (or its uniformity component) (Gupta et al., 2023; Devillers & Lefort, 2022), applying variance regularization (Garrido et al., 2023), or explicitly pushing features apart (van der Pol et al., 2020; Kipf et al., 2020), help mitigate such trivialities, but can inadvertently impair learning of neat equivariance due to their inferior compatibility with equivariance.

Our approach tackles this by directly avoiding the trivial solutions of equivariance:

- **Non-Identity Rotation.** We prevent the trivial case where $\mathcal{R}(\theta_{a(t)})$ becomes the identity by enforcing that the rotation angle $\theta_{a(t)}$ is non-zero. The 2π cyclic case is excluded by design, since θ_a is constrained to $[-\pi, \pi]$. For parameterized augmentations, we introduce a user-defined proxy $\tilde{q}_{a(t)} \in \mathbb{R}^m$, serving as a lower bound on the style variation magnitude, and impose:

$$\mathcal{L}_{\text{theta}} = \mathbb{E}_{a(t) \in \mathcal{A}} \max(0, \epsilon * \tilde{q}_{a(t)} - |\theta_{a(t)}|), \quad (10)$$

where $\theta_{a(t)} \in \mathbb{R}^m$ is the per-plane induced angles by $a(t)$ and $\epsilon > 0$ is a small margin. Note that for subspaces that are unaffected by $a(t)$, $|\theta_{a(t)}|$ is constantly zero by definition.

Table 1: DCI metrics comparison using Random Forest regression on four distinct datasets, shown as triples D/C/I for Disentanglement/Completeness/Informativeness, respectively.

Method/Dataset	MPI3D (D/C/I)	Shape3D (D/C/I)	3DIdent (D/C/I)	3DIEBench (D/C/I)
<i>(Self-Supervised Invariance)</i>				
SimCLR (Chen et al., 2020a)	0.085 / 0.086 / 0.876	0.074 / 0.071 / 0.885	0.129 / 0.096 / 0.926	0.102 / 0.095 / 0.904
<i>(Self-Supervised Equivariance)</i>				
EquiMOD (Devillers & Lefort, 2022)	0.143 / 0.127 / 0.959	0.125 / 0.118 / 0.968	0.067 / 0.024 / 0.933	0.091 / 0.088 / 0.911
CARE(Gupta et al., 2023)	0.084 / 0.085 / 0.891	0.078 / 0.0770 / 0.925	0.147 / 0.119 / 0.928	0.116 / 0.106 / 0.942
<i>(Self-Supervised Disentanglement)</i>				
IP-IRM (Wang et al., 2021)	0.220 / 0.199 / 0.955	0.135 / 0.110 / 0.952	0.147 / 0.094 / 0.933	0.109 / 0.107 / 0.929
Eastwood et al. (2023)	0.950 / 0.826 / 1.000	0.937 / 0.766 / 1.000	0.943 / 0.772 / 0.996	0.362 / 0.297 / 0.952
CD-RED (before post-proc.)	0.628 / 0.505 / 1.000	0.630 / 0.494 / 1.000	0.989 / 0.940 / 1.000	0.595 / 0.481 / 1.000
CD-RED (after post-proc.)	1.000 / 1.000 / 1.000	1.000 / 1.000 / 1.000	0.987 / 0.987 / 1.000	0.996 / 0.996 / 1.000

- **Non-Zero Subspace Norm.** Since rotating a zero vector yields zero, we encourage each 2D plane to maintain a non-trivial radius close to a target $\omega > 0$:

$$\mathcal{L}_{\text{radius}} = \mathbb{E}_{x \in \mathcal{X}} \left[\sum_{i=1}^m \left(\left\| \frac{[H_{c(x)} f(x)]_i}{\omega} \right\| - 1 \right)^2 \right].$$

We require $m\omega^2 < 1$ (with smaller ω preferred) to leave sufficient norm for the non-rotational complement $([H_{c(x)} f(x)]_{\text{extra}})$ to enhance cluster separation.

With the target radius established, we arrive at our equivariance objective. We divide the features by ω as for $\mathcal{L}_{\text{radius}}$ to eliminate sensitivity to the hyperparameter ω :

$$\mathcal{L}_{\text{equi}} = \mathbb{E}_{x \in \mathcal{X}, a(t) \in \mathcal{A}} \left[\left\| \frac{H_{c(x)} f(a(t, x))}{\omega} - \frac{\mathcal{R}(\theta_{a(t)}) H_{c(x)} f(x)}{\omega} \right\|^2 \right]. \quad (11)$$

Our final training objective for CD-RED is thus given by:

$$\mathcal{L}(f, \{g_i\}_i) = \mathcal{L}_{\text{equi}} + \lambda_{\text{radius}} \mathcal{L}_{\text{radius}} + \lambda_{\text{theta}} \mathcal{L}_{\text{theta}}. \quad (12)$$

4 THEORETICAL AND EMPIRICAL RESULTS

In this section, we present theoretical results and empirical validation of the proposed CD-RED from Section 3. First, we show how our framework achieves neat equivariance and subsequent disentanglement, following the definition of Higgins et al. (2018), by leveraging transformations that induce orthogonal variations. Second, we extend the analysis to more complex transformations that may influence multiple latent factors. Formal statements and proofs are provided in Appendix C.

4.1 NEAT EQUIVARIANCE

Proposition 1 (Perfect Equivariance). *CD-RED is capable of achieving perfect equivariance, in the sense that all losses $(\mathcal{L}_{\text{equi}}, \mathcal{L}_{\text{theta}}, \mathcal{L}_{\text{radius}})$ can simultaneously reach their optimum, where each augmentation $a(t)$ that induces a non-trivial variation in the input will correspondingly produce a non-trivial and unique transformation of the feature.*

The idea behind it is simple: by design, all these losses operate on separate aspects of the features and do not contradict each other, unlike most existing solutions (Devillers & Lefort, 2022; Garrido et al., 2023; Gupta et al., 2023).

4.2 AUGMENTATIONS INDUCE ORTHOGONAL VARIATIONS

We now show that achieving perfect equivariance, as defined above, already implies a basic form of disentanglement when augmentations act on distinct subspaces.

Proposition 2 (Weak Disentanglement). *If perfect equivariance is achieved and under mild conditions, for any rotation subspace i assigned exclusively to an orthogonal augmentation $a_i(t)$, the learned angle is linear in the induced latent displacement $q_{a_i}(t)$:*

$$\theta_{a_i(t)} = g_i(t) = c_0^{(i)} q_{a_i}(t), \quad c_0^{(i)} \neq 0,$$

and only the i -th plane changes under $a_i(t)$:

$$[H_{c(x)}f(a(t, x))]_i = \mathcal{R}(\theta_{a_i(t)}^{(i)}) [H_{c(x)}f(x)]_i, \quad [H_{c(x)}f(a(t, x))]_j = [H_{c(x)}f(x)]_j \quad \forall j \neq i.$$

That is, when the augmentations are orthogonal, the rotation parameters in their exclusively assigned subspaces linearly correspond to the latent displacement induced by the augmentations.

A limitation of this weak form is that the feature angle in a subspace may not monotonically track the variation $q_{a(t)}$ due to 2π wrap-around (Figure 2). To prevent ambiguous crossings, we constrain $\theta_{a(t)}$ so that even the corresponding rotation angle for the largest rotations in the data remains within $[-\pi, \pi]$, avoiding wrap-around:

$$\theta_{a(t)} \in \eta \cdot \frac{\max_{t \in T_a} \tilde{q}_{a(t)}}{\max_{t \in T_{\text{data}}} \tilde{q}_a(t)} \cdot [-\pi, \pi], \quad (13)$$

where $\eta \in (0, 1]$ is a usage threshold, e.g., $\eta = 0.5$ for non-cyclic latents and $\eta = 1$ for cyclic. We use $\max_{t \in T_a} \tilde{q}_{a(t)}$ and $\max_{t \in T_{\text{data}}} \tilde{q}_a(t)$ to denote the approximate maximum variation strength present in the augmentation and data, respectively.

Proposition 3 (Strong Disentanglement). *If $\theta_{a(t)}$ is constrained in range to avoid wrap-around as in Eq. (13), then the subspace features associated with $a(t)$ vary **monotonically and linearly** with the underlying style variation $q_{a(t)}$ in the data, achieving strong disentanglement.*

With $\theta_{a(t)}$ effectively constrained, we can further derive a shared, interpretable coordinate through post-processing. For any data $x = \varphi(z)$ and subspace j write

$$[H_{c(x)}f(x)]_j := (u_j(x), v_j(x)), \quad \theta^{(j)}(x) := \text{atan2}(v_j(x), u_j(x)) \in (-\pi, \pi].$$

To obtain an interpretable readout, we align the angles by an offset $\alpha_{c(x),j}$ specific for cluster $c(x)$ and plane j and compute the signed principal difference:

$$\hat{f}^{(j)}(x) := \text{atan2}(\sin(\theta^{(j)}(x) - \alpha_{c(x),j}), \cos(\theta^{(j)}(x) - \alpha_{c(x),j})) \in (-\pi, \pi].$$

- **Noncyclic latents.** When $z^{(j)}$ within a cluster is noncyclic, choose $\alpha_{c(x),j}$ as the midpoint of the covered data arc, then Corollary 2 provides an exact affine readout of $\hat{f}^{(j)}(x)$ relative to latent $z^{(j)}$, and $\hat{f}^{(j)}$ is globally aligned across clusters whenever they share the same latent midpoint.
- **Cyclic latents.** For intrinsically cyclic factors, no lossless 1-D scalarization exists; we therefore retain the 2-D coordinates $[H_{c(x)}f(x)]_j$ on S^1 . By Corollary 3, each cluster c_k differs only by a fixed in-plane rotation $A_{j,k} \in SO(2)$. Global comparison across clusters is achieved via a simple cluster-wise phase alignment, e.g., rotating each cluster by $A_{j,k}^{-1}$ toward a chosen reference.
- **Content latents.** Finally, appending a one-hot cluster indicator to the post-processed style coordinates produces a feature vector disentangled across *style latents* (per subspace as described above) and *content latents* (cluster identity).

Empirical results. We evaluate on **MPI3D** (Gondal et al., 2019) and **Shape3D** (Kim & Mnih, 2018) (primarily discrete latents), as well as on **3DIdent** (Zimmermann et al., 2021) and **3DIEBench** (Garido et al., 2023) (continuous latents). Latent transformations are used as augmentations, with the transformation parameter $t = q_{a(t)}$. We adopt **DCI** (Eastwood & Williams, 2018) as our evaluation metric, which quantifies disentanglement, completeness, and informativeness via supervised regressors mapping features to ground-truth latents. Each component is in the range $[0, 1]$. Larger is

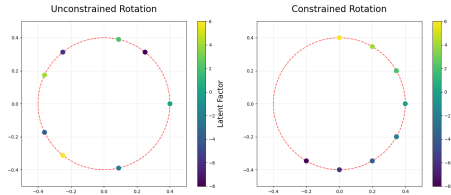


Figure 2: **Left:** Without constraining θ_a , feature subspaces may overlap, causing crossed usage. **Right:** Constraining the range of θ_a prevents overlap, ensuring a monotonic relationship between the variation and the learned feature.

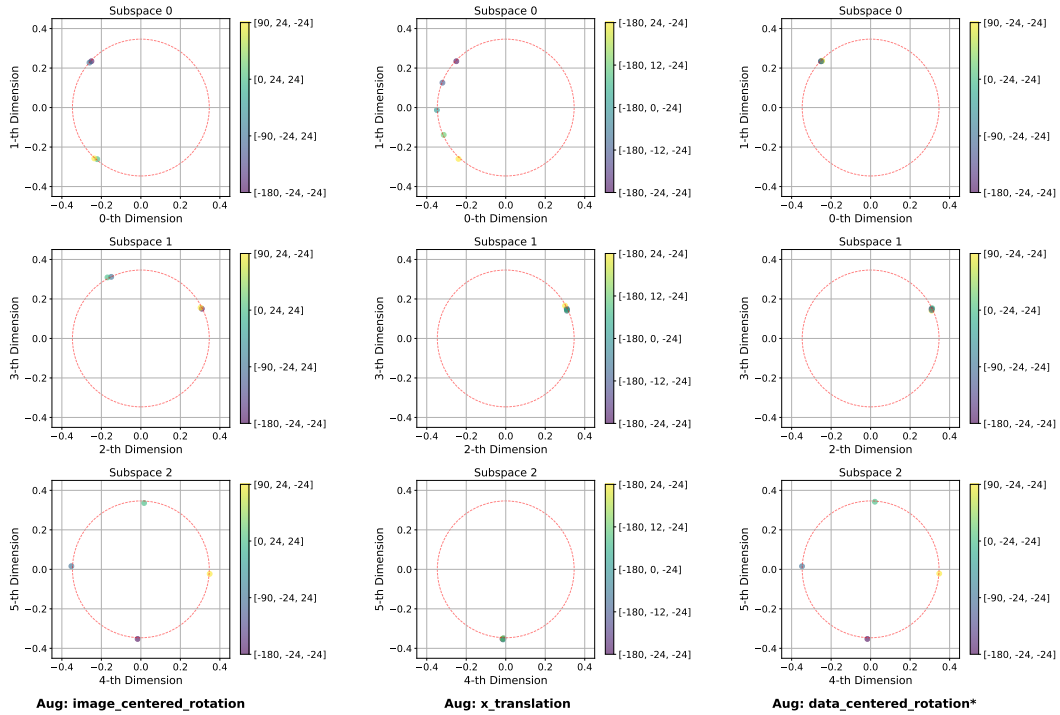


Figure 3: Experiment with x-translation, y-translation, and image-centered rotations. We visualize the learned features in the three subspaces. The image-centered rotations are perfectly disentangled into x-translation, y-translation, and data-centered rotation. The bar indicates the latents of samples.

better. Table 1 compares CD-RED with SimCLR (Chen et al., 2020a), recent self-supervised equivariant baselines (Devillers & Lefort, 2022; Gupta et al., 2023), and self-supervised disentanglement baselines (Wang et al., 2021; Eastwood et al., 2023). Across all four datasets, encompassing both discrete and continuous latent settings, CD-RED achieves near-perfect DCI scores and significantly outperforms the baselines. Details of datasets, metrics, and protocols are provided in Appendix D.

4.3 AUGMENTATIONS INDUCE COMPLEX VARIATIONS

Real-world transformations introduce complex variations that involve multiple underlying factors (e.g., rotating a non-centric object around the origin also induces coupled translations). We now formalize how our framework disentangles such complex augmentations.

Proposition 4 (Disentanglement of Simple Composite Augmentations). *Given m orthogonal augmentations, each assigned exclusively to a unique subspace, if perfect equivariance is achieved, then for any composite augmentation a_{comp} that induces only n ($\leq m$) of these style variations, the framework will learn to disentangle them into the corresponding n orthogonal subspaces, with each variation captured by its associated rotation angle θ .*

This is because the m orthogonal augmentations define how each input variation corresponds to the rotation angles θ in their respective subspaces. Any additional augmentation that overlaps with these subspaces must follow these predefined mappings.

Proposition 5 (Disentanglement of Complex Augmentation with Anchor). *Given $m - 1$ orthogonal augmentations, each assigned to a unique rotation plane, let a_{comp} be an additional augmentation such that there exists an anchor where a_{comp} acts purely on the m -th style latent, and away from the anchor, a_{comp} decomposes into a composite over n ($\leq m$) already defined subspaces. Then, assigning a new dedicated plane m to this anchor yields a linear readout $\theta_{a_{comp}(t)}^{(m)} = c_0^{(m)} q_{comp}(t)$, and subsequently the angle readout in all affected n ($\leq m$) subspaces is additive across planes, as described in Proposition 4.*

This is because the existence of an anchor defines the m -th subspace, and the remaining then directly follows from Proposition 4.

Empirical results. We further provide empirical validation for Proposition 5 in Figure 3. A similar validation of Proposition 4 is provided in Figure 13 in the Appendix F.2.

5 ADDITIONAL RELATED WORK

Self-supervised Clustering. Self-supervised clustering methods are commonly built on either contrastive or non-contrastive objectives. Most contrastive methods (Yeh et al., 2022; He et al., 2020; Chen et al., 2020b;a) use an InfoNCE-type loss (Oord et al., 2018). HaoChen et al. (2021) studies a spectral contrastive loss, and Caron et al. (2018; 2020); Li et al. (2020) treat cluster centroids (prototypes) as contrastive views instead of individual samples. Non-contrastive methods remove the need for explicit negative pairs, either by architectural design (Grill et al., 2020; Chen & He, 2021) or by regularizing the covariance of the embeddings (Ermolov et al., 2021; Zbontar et al., 2021; Bardes et al., 2021; 2022). Although there are works showing that contrastive and non-contrastive methods share theoretical connections (Garrido et al., 2022), Pokle et al. (2022) argue that non-contrastive methods may fail to learn cluster-aligned representations. Our first stage of training builds upon the SimCLR framework (Chen et al., 2020a), for which recent analyses provide theoretical guarantees that the learned features form well-separated clusters in the hypersphere (Wang & Isola, 2020; Wang et al., 2022; Parulekar et al., 2023).

Group-equivariance Representation Learning. Group-equivariance representation learning aims to learn feature maps whose response to structured transformations of the input follows a group representation, i.e., $\varphi(T_g(x)) = \rho(g)\varphi(x)$ for a group $g \in G$. Classical group-equivariance CNNs instantiate this idea by fixing both the group G and its action on the input domain, then designing convolutional layers so that feature channels transform according to a prescribed representation ρ (Cohen & Welling, 2016a;b; Cohen et al., 2018; Weiler et al., 2018; Weiler & Cesa, 2019; Finzi et al., 2020; Lengyel et al., 2023; Yang et al., 2024). More recently, the architecture extends to message-passing GNN (Anderson et al., 2019; Satorras et al., 2021) and Transformer (Tai et al., 2019; Fuchs et al., 2020). However, in these approaches, the group and its feature-space representation are fixed a priori and hard-coded into the architecture. In contrast, our method is more flexible by keeping the backbone architecture generic and instead learning equivariance in feature space.

Disentanglement in Generative Models. Separating latent factors of variation is one of the central objectives in generative modeling, and has been explored in VAE-based approaches (Higgins et al., 2017; Kim & Mnih, 2018; Chen et al., 2018; Li et al., 2019; Mathieu et al., 2019; Esmaili et al., 2019), GAN-based methods (Chen et al., 2016; Jeon et al., 2021; Ojha et al., 2020), and more recent denoising-based generative models (Yang et al., 2023; Song et al., 2023; Wu & Zheng, 2024; Hu et al., 2024). These methods aim to map data to structured latent spaces, but in practice often struggle to balance sample quality, degree of disentanglement, and robustness to complex data distributions. There also exists a line of work that promotes disentanglement through weak supervision in the form of pairwise or grouped observations, which can be viewed as imposing conditional invariance constraints on subsets of the latent variables (Bouchacourt et al., 2018; Hosoya, 2018; Shu et al., 2019; Chen & Batmanghelich, 2020; Locatello et al., 2020). In contrast, our method is purely self-supervised while achieving cleaner and stronger disentanglement.

6 CONCLUSION

In this work, we have proposed CD-RED, a two-stage self-supervised disentanglement framework that learns initial semantic clustering via InfoNCE in the first stage and then learns cluster-dependent rotational equivariance in the second stage. CD-RED aligns rotation coordinates of different clusters via Householder transformations and models equivariance as $SO(2)$ groups within orthogonal hyperspherical subspaces. We have theoretically established that, under mild assumptions, the resulting representations achieve neat equivariance and perfect disentanglement of cluster and style latents. Our results go beyond standard group-disentanglement settings, providing disentanglement under more practical, non-orthogonal augmentations. Extensive empirical results have demonstrated that CD-RED can robustly achieve near-perfect disentanglement across various settings.

REFERENCES

- 540
541
542 Brandon Anderson, Truong Son Hy, and Risi Kondor. Cormorant: Covariant molecular neural
543 networks. *Advances in neural information processing systems*, 32, 2019.
- 544
545 Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization
546 for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- 547
548 Adrien Bardes, Jean Ponce, and Yann LeCun. Vicregl: Self-supervised learning of local visual
549 features. *Advances in Neural Information Processing Systems*, 35:8799–8810, 2022.
- 550
551 Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new
552 perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828,
553 2013.
- 554
555 Diane Bouchacourt, Ryota Tomioka, and Sebastian Nowozin. Multi-level variational autoencoder:
556 Learning disentangled representations from grouped observations. In *Proceedings of the AAAI
557 Conference on Artificial Intelligence*, volume 32, 2018.
- 558
559 Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 560
561 Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for un-
562 supervised learning of visual features. In *Proceedings of the European conference on computer
563 vision (ECCV)*, pp. 132–149, 2018.
- 564
565 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
566 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural
567 information processing systems*, 33:9912–9924, 2020.
- 568
569 Junxiang Chen and Kayhan Batmanghelich. Weakly supervised disentanglement by pairwise simi-
570 larities. *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 3495–3502, Jun 2020.
- 571
572 Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disen-
573 tanglement in variational autoencoders. *Advances in neural information processing systems*, 31,
574 2018.
- 575
576 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
577 contrastive learning of visual representations. In *International conference on machine learning*,
578 pp. 1597–1607. PMLR, 2020a.
- 579
580 Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Info-
581 gan: Interpretable representation learning by information maximizing generative adversarial nets.
582 *Advances in neural information processing systems*, 29, 2016.
- 583
584 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of
585 the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
- 586
587 Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum
588 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- 589
590 Brian Cheung, Jesse A Livezey, Arjun K Bansal, and Bruno A Olshausen. Discovering hidden
591 factors of variation in deep networks. *arXiv preprint arXiv:1412.6583*, 2014.
- 592
593 Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International confer-
ence on machine learning*, pp. 2990–2999. PMLR, 2016a.
- Taco S Cohen and Max Welling. Steerable cnns. *arXiv preprint arXiv:1612.08498*, 2016b.
- Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical cnns. *arXiv preprint
arXiv:1801.10130*, 2018.
- Rumen Dangovski, Li Jing, Charlotte Loh, Seungwook Han, Akash Srivastava, Brian Cheung, Pulkit
Agrawal, and Marin Soljačić. Equivariant contrastive learning. *arXiv preprint arXiv:2111.00899*,
2021.

- 594 Alexandre Devillers and Mathieu Lefort. Equimod: An equivariance module to improve self-
595 supervised learning. *arXiv preprint arXiv:2211.01244*, 2022.
596
- 597 Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of disen-
598 tangled representations. In *6th International Conference on Learning Representations*, 2018.
- 599 Cian Eastwood, Julius von Kügelgen, Linus Ericsson, Diane Bouchacourt, Pascal Vincent, Bernhard
600 Schölkopf, and Mark Ibrahim. Self-supervised disentanglement by leveraging structure in data
601 augmentations. *arXiv preprint arXiv:2311.08815*, 2023.
602
- 603 Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-
604 supervised representation learning. In *International conference on machine learning*, pp. 3015–
605 3024. PMLR, 2021.
- 606 Babak Esmaeili, Hao Wu, Sarthak Jain, Alican Bozkurt, Narayanaswamy Siddharth, Brooks Paige,
607 Dana H Brooks, Jennifer Dy, and Jan-Willem Meent. Structured disentangled representations.
608 In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2525–2534.
609 PMLR, 2019.
- 610 Marc Finzi, Samuel Stanton, Pavel Izmailov, and Andrew Gordon Wilson. Generalizing convolu-
611 tional neural networks for equivariance to lie groups on arbitrary continuous data. In *International
612 conference on machine learning*, pp. 3165–3176. PMLR, 2020.
613
- 614 Fabian Fuchs, Daniel Worrall, Volker Fischer, and Max Welling. Se (3)-transformers: 3d roto-
615 translation equivariant attention networks. *Advances in neural information processing systems*,
616 33:1970–1981, 2020.
- 617 Quentin Garrido, Yubei Chen, Adrien Bardes, Laurent Najman, and Yann Lecun. On the duality be-
618 tween contrastive and non-contrastive self-supervised learning. *arXiv preprint arXiv:2206.02574*,
619 2022.
620
- 621 Quentin Garrido, Laurent Najman, and Yann Lecun. Self-supervised learning of split invariant
622 equivariant representations. *arXiv preprint arXiv:2302.10283*, 2023.
- 623 Muhammad Waleed Gondal, Manuel Wuthrich, Djordje Miladinovic, Francesco Locatello, Martin
624 Breidt, Valentin Volchkov, Joel Akpo, Olivier Bachem, Bernhard Schölkopf, and Stefan Bauer.
625 On the transfer of inductive bias from simulation to the real world: a new disentanglement dataset.
626 *Advances in Neural Information Processing Systems*, 32, 2019.
627
- 628 Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre Richemond, Elena
629 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,
630 et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural
631 information processing systems*, 33:21271–21284, 2020.
- 632 Sharut Gupta, Joshua Robinson, Derek Lim, Soledad Villar, and Stefanie Jegelka. Structur-
633 ing representation geometry with rotationally equivariant contrastive learning. *arXiv preprint
634 arXiv:2306.13924*, 2023.
- 635
- 636 Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-supervised
637 deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*,
638 34:5000–5011, 2021.
- 639 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
640 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on
641 computer vision and pattern recognition*, pp. 9729–9738, 2020.
- 642 Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew M Botvinick,
643 Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a
644 constrained variational framework. *ICLR (Poster)*, 3, 2017.
645
- 646 Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende,
647 and Alexander Lerchner. Towards a definition of disentangled representations. *arXiv preprint
arXiv:1812.02230*, 2018.

- 648 Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural
649 networks. *science*, 313(5786):504–507, 2006.
- 650 Kurt Hornik, Ingo Feinerer, Martin Kober, and Christian Buchta. Spherical k-means clustering.
651 *Journal of statistical software*, 50:1–22, 2012.
- 652 Haruo Hosoya. Group-based learning of disentangled representations with generalizability for novel
653 contents. *arXiv preprint arXiv:1809.02383*, 2018.
- 654 Alston S. Householder. Unitary triangularization of a nonsymmetric matrix. *Numerische Mathe-*
655 *matik*, 2(2):163–168, 1958.
- 656 Vincent Tao Hu, Wei Zhang, Meng Tang, Pascal Mettes, Deli Zhao, and Cees Snoek. Latent space
657 editing in transformer-based flow matching. In *Proceedings of the AAAI conference on artificial*
658 *intelligence*, volume 38, pp. 2247–2255, 2024.
- 659 Insu Jeon, Wonkwang Lee, Myeongjang Pyeon, and Gunhee Kim. Ib-gan: Disentangled representa-
660 tion learning with information bottleneck generative adversarial networks. In *Proceedings of the*
661 *AAAI conference on artificial intelligence*, volume 35, pp. 7926–7934, 2021.
- 662 Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International conference on ma-*
663 *chine learning*, pp. 2649–2658. PMLR, 2018.
- 664 Thomas Kipf, Elise van der Pol, and Max Welling. Contrastive learning of structured world models.
665 In *International Conference on Learning Representations*, 2020.
- 666 Attila Lengyel, Ombretta Strafforello, Robert-Jan Brintjes, Alexander Gielisse, and Jan van
667 Gemert. Color equivariant convolutional networks. *Advances in Neural Information Process-*
668 *ing Systems*, 36:29831–29850, 2023.
- 669 Junnan Li, Pan Zhou, Caiming Xiong, and Steven CH Hoi. Prototypical contrastive learning of
670 unsupervised representations. *arXiv preprint arXiv:2005.04966*, 2020.
- 671 Yang Li, Quan Pan, Suhang Wang, Haiyun Peng, Tao Yang, and Erik Cambria. Disentangled varia-
672 tional auto-encoder for semi-supervised learning. *Information Sciences*, 482:73–85, 2019.
- 673 Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard
674 Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning
675 of disentangled representations. In *international conference on machine learning*, pp. 4114–4124.
676 PMLR, 2019.
- 677 Francesco Locatello, Ben Poole, Gunnar Raetsch, Bernhard Schölkopf, Olivier Bachem, and
678 Michael Tschannen. Weakly-supervised disentanglement without compromises. *International*
679 *Conference on Machine Learning, International Conference on Machine Learning*, 2020.
- 680 Giovanni Luca Marchetti, Gustaf Tegnér, Anastasiia Varava, and Danica Kragic. Equivariant repre-
681 sentation learning via class-pose decomposition. In *International Conference on Artificial Intel-*
682 *ligence and Statistics*, pp. 4745–4756. PMLR, 2023.
- 683 Emile Mathieu, Tom Rainforth, Nana Siddharth, and Yee Whye Teh. Disentangling disentanglement
684 in variational autoencoders. In *International conference on machine learning*, pp. 4402–4412.
685 PMLR, 2019.
- 686 Daniel Müllner. Modern hierarchical, agglomerative clustering algorithms. *arXiv preprint*
687 *arXiv:1109.2378*, 2011.
- 688 Lilian Ngweta, Subha Maity, Alex Gittens, Yuekai Sun, and Mikhail Yurochkin. Simple disentan-
689 glement of style and content in visual representations. In *International Conference on Machine*
690 *Learning*, pp. 26063–26086. PMLR, 2023.
- 691 Utkarsh Ojha, Krishna Kumar Singh, Cho-Jui Hsieh, and Yong Jae Lee. Elastic-infogan: Unsuper-
692 vised disentangled representation learning in class-imbalanced data. *Advances in neural informa-*
693 *tion processing systems*, 33:18063–18075, 2020.

- 702 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predic-
703 tive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 704
- 705 Advait Parulekar, Liam Collins, Karthikeyan Shanmugam, Aryan Mokhtari, and Sanjay Shakkot-
706 tai. Infonce loss provably learns cluster-preserving representations. In *The Thirty Sixth Annual*
707 *Conference on Learning Theory*, pp. 1914–1961. PMLR, 2023.
- 708 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
709 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
710 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- 711 Ashwini Pokle, Jinjin Tian, Yuchen Li, and Andrej Risteski. Contrasting the landscape of contrastive
712 and non-contrastive learning. *arXiv preprint arXiv:2203.15702*, 2022.
- 713
- 714 Robin Quessard, Thomas D. Barrett, and William R. Clements. Learning group structure and disen-
715 tangled representations of dynamical environments, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2002.06991)
716 [2002.06991](https://arxiv.org/abs/2002.06991).
- 717 Scott Reed, Kihyuk Sohn, Yuting Zhang, and Honglak Lee. Learning to disentangle factors of
718 variation with manifold interaction. In *International conference on machine learning*. PMLR,
719 2014.
- 720 Xuanchi Ren, Tao Yang, Yuwang Wang, and Wenjun Zeng. Rethinking content and style: exploring
721 bias for unsupervised disentanglement. In *Proceedings of the IEEE/CVF International Confer-*
722 *ence on Computer Vision*, pp. 1823–1832, 2021.
- 723
- 724 Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E (n) equivariant graph neural net-
725 works. In *International conference on machine learning*, pp. 9323–9332. PMLR, 2021.
- 726 Mehran Shakerinava, Arnab Kumar Mondal, and Siamak Ravanbakhsh. Structuring representations
727 using group invariants. *Advances in Neural Information Processing Systems*, 35:34162–34174,
728 2022.
- 729
- 730 Rui Shu, Yining Chen, Abhishek Kumar, Stefano Ermon, and Ben Poole. Weakly supervised disen-
731 tanglement with guarantees. *arXiv preprint arXiv:1910.09772*, 2019.
- 732 Yue Song, Andy Keller, Nicu Sebe, and Max Welling. Flow factorized representation learning.
733 *Advances in Neural Information Processing Systems*, 36:49761–49782, 2023.
- 734 Kai Sheng Tai, Peter Bailis, and Gregory Valiant. Equivariant transformer networks. In *International*
735 *Conference on Machine Learning*, pp. 6086–6095. PMLR, 2019.
- 736
- 737 Pieter Merkus Lambertus Tammes. On the origin of number and arrangement of the places of exit
738 on the surface of pollen-grains. *Recueil des travaux botaniques néerlandais*, 27(1):1–84, 1930.
- 739
- 740 Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Computer*
741 *Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings,*
742 *Part XI 16*, pp. 776–794. Springer, 2020.
- 743 Elise van der Pol, Thomas Kipf, Frans A Oliehoek, and Max Welling. Plannable approximations
744 to mdp homomorphisms: Equivariance under actions. In *Proceedings of the 19th International*
745 *Conference on Autonomous Agents and MultiAgent Systems*, pp. 1431–1439, 2020.
- 746 Julius Von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel
747 Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably
748 isolates content from style. *Advances in neural information processing systems*, 34:16451–16467,
749 2021.
- 750
- 751 Tan Wang, Zhongqi Yue, Jianqiang Huang, Qianru Sun, and Hanwang Zhang. Self-supervised
752 learning disentangled group representation as feature. *Advances in Neural Information Processing*
753 *Systems*, 34:18225–18240, 2021.
- 754
- 755 Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through align-
ment and uniformity on the hypersphere. In *International conference on machine learning*, pp.
9929–9939. PMLR, 2020.

- 756 Yifei Wang, Qi Zhang, Yisen Wang, Jiansheng Yang, and Zhouchen Lin. Chaos is a ladder: A
757 new theoretical understanding of contrastive learning via augmentation overlap. *arXiv preprint*
758 *arXiv:2203.13457*, 2022.
- 759 Maurice Weiler and Gabriele Cesa. General $e(2)$ -equivariant steerable cnns. *Advances in neural*
760 *information processing systems*, 32, 2019.
- 761 Maurice Weiler, Mario Geiger, Max Welling, Wouter Boomsma, and Taco S Cohen. 3d steerable
762 cnns: Learning rotationally equivariant features in volumetric data. *Advances in Neural informa-*
763 *tion processing systems*, 31, 2018.
- 764 Ancong Wu and Wei-Shi Zheng. Factorized diffusion autoencoder for unsupervised disentangled
765 representation learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol-
766 *ume 38*, pp. 5930–5939, 2024.
- 767 Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-
768 parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision*
769 *and pattern recognition*, pp. 3733–3742, 2018.
- 770 Tete Xiao, Xiaolong Wang, Alexei A Efros, and Trevor Darrell. What should not be contrastive in
771 contrastive learning. *arXiv preprint arXiv:2008.05659*, 2020.
- 772 Tao Yang, Yuwang Wang, Yan Lv, and Nanning Zheng. Disdiff: Unsupervised disentanglement of
773 diffusion probabilistic models. *arXiv preprint arXiv:2301.13721*, 2023.
- 774 Yulong Yang, Felix O’Mahony, and Christine Allen-Blanchette. Learning color equivariant repre-
775 sentations. *arXiv preprint arXiv:2406.09588*, 2024.
- 776 Chun-Hsiao Yeh, Cheng-Yao Hong, Yen-Chi Hsu, Tyng-Luh Liu, Yubei Chen, and Yann LeCun. De-
777 coupled contrastive learning. In *European conference on computer vision*, pp. 668–684. Springer,
778 2022.
- 779 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised
780 learning via redundancy reduction. In *International conference on machine learning*, pp. 12310–
781 12320. PMLR, 2021.
- 782 Roland S Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel.
783 Contrastive learning inverts the data generating process. In *International Conference on Machine*
784 *Learning*, pp. 12979–12990. PMLR, 2021.
- 785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A USE OF LARGE LANGUAGE MODELS

We used a large language model (LLM) only for language polishing (grammar, wording, and clarity) of drafts written by the authors. The model did not generate research ideas, methods, analyses, results, or figures, and it did not write any sections from scratch.

B EXTENSIONS AND SETTINGS OF FIGURE 1

Settings. We construct a toy dataset of three uppercase letters (“R”, “E”, “D”) rendered as 96×96 grayscale images. Each instance is generated by rotating the glyph over angles $\{-165^\circ, -150^\circ, \dots, 180^\circ\}$ (step 15°), forming a full cyclic orbit. During training, we use only a 15° rotation augmentation (bilinear interpolation). All models use temperature $\tau = 0.5$.

We follow Wang & Isola (2020) to decompose the contrastive objective InfoNCE into *invariance* and *uniformity* terms. For CARE (Gupta et al., 2023), we sweep the weights on these three components (invariance/uniformity/equivariance) to produce the results in Figure 1.

Empirically, the three terms *compete*: stronger equivariance encourages features to align along rotation orbits, whereas strong invariance/uniformity tends to collapse or spread them irrespective of orbit structure. The best equivariance is achieved when cluster layouts are aligned with the rotational orbits (i.e., parallel).

Extension: non-cyclic setting. We repeat the visualization on a non-cyclic range by restricting angles to $[-120^\circ, 120^\circ]$ (same dataset, step 15°).

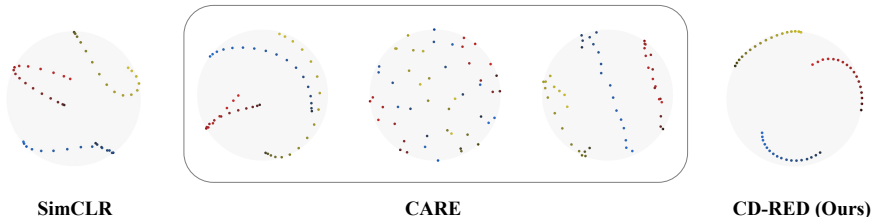


Figure 4: Comparisons for non-cyclic setting.

C PROOFS AND REMARKS

C.1 PRELIMINARY

Definition 1 (Augmentation Connectivity). Let \mathcal{A} be a set of augmentations acting on the observation space \mathcal{X} . Define a relation $\sim_{\mathcal{A}}$ on \mathcal{X} by

$$x \sim_{\mathcal{A}} y \iff \exists a, b \in \mathcal{A} : \varphi^{-1}(a(x)) = \varphi^{-1}(b(y)).$$

The relation $\sim_{\mathcal{A}}$ is an equivalence relation:

1. **Reflexive.** $x \sim_{\mathcal{A}} x$;
2. **Symmetric.** If $x \sim_{\mathcal{A}} y$, then $y \sim_{\mathcal{A}} x$;
3. **Transitive.** If $x \sim_{\mathcal{A}} y$, $y \sim_{\mathcal{A}} z$, then $x \sim_{\mathcal{A}} z$.

A subset $D \subseteq \mathcal{X}$ is said to be \mathcal{A} -connected if $x \sim_{\mathcal{A}} y$ for every pair $x, y \in D$.

Assumption 1 (Intra-cluster Connectivity). Fix a collection of augmentations \mathcal{A} acting on the observation space \mathcal{X} and let $\sim_{\mathcal{A}}$ denote the augmentation-connectivity relation defined in Definition 1. Let $\mathcal{X}_{\text{train}} \subset \mathcal{X}$ be the training set, and write

$$\mathcal{Z}_{\text{train}} := \varphi^{-1}(\mathcal{X}_{\text{train}}) \subseteq U_{\mathcal{A}} \times \mathcal{S}_{\mathcal{A}},$$

where $U_{\mathcal{A}}$ is the cluster subspace and $\mathcal{S}_{\mathcal{A}}$ is the style subspace.

864 **1. Augmentation-wise factorization.** Each $a \in \mathcal{A}$ lifts to a map

$$865 \tilde{a} : \mathcal{S}_{\mathcal{A}} \times \mathcal{U}_{\mathcal{A}} \longrightarrow \mathcal{S}_{\mathcal{A}} \times \mathcal{U}_{\mathcal{A}}, \quad \tilde{a}(s, u) = (\tilde{s}, u),$$

866 where $\tilde{s} \neq s$ for non-identity transformations. Thus, every augmentation acts only on the
867 style coordinates s and leaves the cluster coordinates u unchanged.
868

869 **2. Intra-cluster \mathcal{A} -connectivity.** Define the cluster partitions

$$870 \mathcal{C}_{\mathcal{A}} := \{c_u := \{x \in \mathcal{X}_{\text{train}} : \varphi^{-1}(x) = (s, u), s \in \mathcal{S}_{\mathcal{A}}\}, \forall u \in \mathcal{U}_{\mathcal{A}}\}.$$

871 We assume that each c_u is \mathcal{A} -connected: for all $x, y \in c_u$ one has $x \sim_{\mathcal{A}} y$.
872

873 Because augmentations never modify the u -coordinate, no sequence of augmentations can connect
874 points with different cluster coordinates, i.e., if $x \in c_u$ and $y \in c_{u'}$ with $u \neq u'$ then $x \not\sim_{\mathcal{A}} y$.

875 **Definition 2** (Cluster-wise axis-aligned representation). Let $f : \mathcal{X} \rightarrow \mathbb{S}^{d-1}$ be a unit-norm rep-
876 resentation. Fix an augmentation set \mathcal{A} that alters exactly d_{sty} style factors and induces cluster
877 partitions $\mathcal{C}_{\mathcal{A}}$ of $\mathcal{X}_{\text{train}}$. We say f is **cluster-wise axis-aligned** if there exist orthogonal matrices
878 $\{H_c \in O(n)\}_{c \in \mathcal{C}_{\mathcal{A}}}$ (one per cluster) such that, for the aligned representation

$$879 \tilde{f}(x) := H_{c(x)} f(x) = (\tilde{f}_{\text{sty}}^{(1)}(x), \dots, \tilde{f}_{\text{sty}}^{(d_{\text{sty}})}(x), \tilde{f}_{d-2d_{\text{sty}}-1}(x)) \in (\mathbb{S}^1)^{d_{\text{sty}}} \times \mathbb{S}^{d-2d_{\text{sty}}-1},$$

880 The following holds:
881

882 **1. Center synchronization.** For each cluster c , H_c maps its rotation center r_c to the north
883 pole $e_d = (0, \dots, 0, 1)^T$ (i.e., $H_c r_c = e_d$).
884

885 **2. Style-plane separation.** For each $1 \leq i \leq d_{\text{sty}}$, the component $\tilde{f}_{\text{sty}}^{(i)}(x) \in \mathbb{S}^1$ lies in the
886 i -th canonical rotation plane $\pi_i := \text{span}\{e_{2i-1}, e_{2i}\}$. Equivalently, for any angle $\theta \in \mathbb{R}$,
887 define the block-diagonal rotation

$$888 \mathcal{R}_i(\theta) := \text{diag}(I_2, \dots, \underbrace{\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}}_{i\text{-th block}}, \dots, I_2, I_{n-2d_{\text{sty}}}).$$

889 Then $\mathcal{R}_i(\theta)$ rotates only $\tilde{f}_{\text{sty}}^{(i)}$ and leaves $\tilde{f}_{\text{sty}}^{(j)}$ ($j \neq i$) and \tilde{f}_{res} unchanged.
890

891 Consequently, we focus solely on disentangling the augmentation-induced style latents in the aligned
892 representation.
893

894 **Definition 3** (Augmentation-Induced Style Latent Variation). Let $a(t) \in \mathcal{A}$ be an augmentation
895 parameterized by $t \in \mathcal{T}_a$, its induced style displacement at latent point $z \in \mathcal{Z}_{\text{train}}$ is defined as:
896

$$897 \Delta_a(t, z) := \varphi^{-1}(a(t, \varphi(z))) - z \in \mathbb{R}^m,$$

898 which is continuous and bijective in t for fixed z .
899

900 In general, the variation $\Delta_a(t, z)$ may depend on both the augmentation parameter t and the latent
901 state z . For example, a centric rotation may induce different latent changes at different positions of
902 the object.

903 **Definition 4** (\mathcal{A} -realized cyclic vs. non-cyclic on plane j). Let $z^{(j)}$ be the j -th style latent and let
904 a_j denote the augmentation varying only this coordinate.
905

906 **1. Intrinsic cyclicity.** We say $z^{(j)}$ is intrinsically cyclic if it has a latent period $p > 0$ (i.e.
907 $z^{(j)} \equiv z^{(j)} + p$, so values live on a circle). Otherwise, it is intrinsically non-cyclic (values
908 live on an interval).
909

910 **2. \mathcal{A} -realized cyclicity.** On the training data support, we say the j -factor is **\mathcal{A} -realized cyclic**
911 if:

912 (a) $z^{(j)}$ is intrinsically cyclic; and

913 (b) there exists an ordering of points (x_1, x_2, \dots, x_K) varying only in $z^{(j)}$ such that

$$914 a_i(t_k, x_k) = a_i(t_{k+1}, x_{k+1}) \quad (k = 1, \dots, K-1), \quad a_i(t_\ell, x_k) = a_i(t_m, x_1),$$

915 Intuitively, the augmentation steps can be chained to form a circle through each sample.
916

917 **3. \mathcal{A} -non-cyclic.** If either the latent is intrinsically non-cyclic, or no such nontrivial closing
chain exists, then the j -factor is \mathcal{A} -non-cyclic.

C.2 PROOF OF PROPOSITION 1

Setup. Stage 1 pretraining yields unit features $f(x)$ with k spherical k -means centers $\{r_\ell\}_{\ell=1}^k \subset \mathbb{S}^{d-1}$. In Stage 2, we freeze the cluster assignment rule $c(x) = \arg \min_\ell \|u(x) - r_\ell\|_2$, build per-cluster Householder maps H_ℓ that send r_ℓ to the north pole, and optimize on losses $\mathcal{L}_{\text{equiv}}, \mathcal{L}_{\text{radius}}, \mathcal{L}_{\text{theta}}$.

Lemma 1. *Under the Assumption 1, if we have a relatively small radius ω such that $d_{\min} := \min_{\ell \neq r} \|c_\ell - c_r\|_2 \leq 2\omega\sqrt{m}$ and $\mathcal{L}_{\text{radius}} \rightarrow 0$, then for any x in cluster ℓ and any augmentation $a(t)$, both x and $a(t, x)$ remain assigned to ℓ .*

Proof. Since all intra-cluster data are \mathcal{A} -connected, in main content, we know that all the clusters can be separated uniformly in different region after pretraining with Infonce loss, so as to the cluster centers, so there is possible that we found enough small radius ω satisfying the requirements.

When $\mathcal{L}_{\text{radius}} = 0$, both $H_\ell f(x)$ and $H_\ell f(a(t, x))$ lie on the product torus $S^1(\omega)^{\times m}$ in the aligned chart, so their *pairwise* chord distance is at most $2\omega\sqrt{m}$ (sum of m plane-wise chord bounds). Hence

$$\|f(a(t, x)) - r_\ell\| \leq \|f(a(t, x)) - f(x)\| + \|f(x) - r_\ell\| \leq 2\omega\sqrt{m} + \|f(x) - r_\ell\|.$$

If some other center r_o were closer to $f(a(t, x))$, then by the triangle inequality

$$\|r_\ell - r_o\| \leq \|r_\ell - f(a(t, x))\| + \|f(a(t, x)) - r_o\| < 2\omega\sqrt{m},$$

contradicting $d_{\min} > 2\omega\sqrt{m}$. Thus $c(a(t, x)) = \ell$. \square

Proposition 6 (Convergence to perfect equivariance with three losses). *Under Lemma. 1, consider Stage 2 with losses $\mathcal{L}_{\text{equiv}}, \mathcal{L}_{\text{radius}}, \mathcal{L}_{\text{theta}}$. There exists a solution at which all three reach their optima simultaneously, and all data x and augmentation a ,*

$$H_{c(x)} f(a(x, t)) = \mathcal{R}(\theta_{a(t)}) H_{c(x)} f(x) \quad \text{for all } x \sim \mathcal{X}, a(t) \in \mathcal{A},$$

Moreover, any nontrivial a induces a non-identity $\mathcal{R}(\theta_{a(t)})$.

Proof. (i) *Fixed cluster charts.* By Lemma 1, assignments are stable, so a single Householder chart H_ℓ consistently aligns both x and $a(t, x)$ to the same “north-pole” frame.

(ii) *Radial normalization.* $\mathcal{L}_{\text{radius}} = 0$ puts features on the product of circles with radius ω in each style plane. Thus aligned features are points on $S^1(\omega)^{\times m}$, and any valid transformation within the representation must be a block-diagonal rotation $\in SO(2)^m$.

(iii) *Non-degenerate angles.* By $\mathcal{L}_{\text{theta}}$, the angle predictor G is constrained so that $\theta_{a(t)}$ varies with t and is bounded to avoid wrap-through ambiguities; in particular, nontrivial t cannot map to $\theta_{a(t)} = 0$.

(iv) *Equivariance identification.* Minimizing $\mathcal{L}_{\text{equiv}}$ forces $\mathcal{R}(\theta_{a(t)}) H_{c(x)} f(x) = H_{c(x)} f(a(t, x))$ point-wise. On $S^1(\omega)^{\times m}$ this equality pins down $\mathcal{R}(\theta_{a(t)})$ uniquely per t (up to 2π per axis, ruled out by the range constraint in $\mathcal{L}_{\text{theta}}$). Since nontrivial t produce $\mathcal{R}(\theta_{a(t)}) H_{c(x)} f(x) = H_{c(x)} f(a(t, x))$, the corresponding $\mathcal{R}(\theta_{a(t)})$ is non-identity. \square

Remark 1 (Robustness to imperfect pretraining clusters). *Even if Stage 1 yields some misassigned samples, Stage 2 can recover from it. Under Assumption 1 (intra-cluster \mathcal{A} -connectivity) and a small radius ω with margin $d_{\min} > 2\omega\sqrt{m}$, the combined losses $\mathcal{L}_{\text{equiv}} + \lambda_{\text{radius}} \mathcal{L}_{\text{radius}} + \lambda_{\text{theta}} \mathcal{L}_{\text{theta}}$ induce the following behavior:*

- **Attraction to the correct chart.** *For a misassigned x , the equivariance residual is strictly smaller in its true cluster chart than in its current one (the wrong Householder frame breaks torus consistency). Gradients of $\mathcal{L}_{\text{equiv}}$ therefore move $u(x)$ toward the Voronoi cell of the correct center.*
- **Bounded motion & stability once corrected.** *$\mathcal{L}_{\text{radius}}$ confines in-plane motion to the torus $S^1(\omega)^{\times m}$, so the maximal within-cluster displacement during augmentation is $\leq 2\omega\sqrt{m}$. Once $f(x)$ crosses into the correct Voronoi cell, Lemma 1 applies and the assignment remains stable thereafter.*

Now, we consider the case that $\Delta_a(t, z)$ is independent of z , i.e., the induced variation depends only on the augmentation parameter t , never on the current latent state.

972 C.3 PROOF OF PROPOSITION 2

973
974 **Lemma 2.** Let $\mathcal{A} = \{a_i(t_i)\}_{i=1}^{d_{\text{sty}}}$ be d_{sty} parametric augmentations, each acting solely on, without
975 loss of generality, the i -th style coordinate and assigning to i th canonical rotation plane. For every
976 latent $z \in \mathcal{Z}_{\text{train}}$ and parameter $t_i \in \mathcal{T}_{a_i}$, define

$$977 \Delta_{a_i}(t_i, z) = q_{a_i}(t_i) e_i, \text{ where } |q_{a_i}(t_i)| \leq M_i q_0^{(i)},$$

978 where e_i is the i -th standard basis vector of \mathbb{R}^m , $q_0^{(i)} > 0$ the canonical step and $M_i \in \mathbb{N}_+$ a
979 coordinate bound. For $k \in \mathbb{N}^{d_{\text{sty}}}$ abbreviate $k \odot q_0 := \sum_j k^{(j)} q_0^{(j)} e_j$.

982 A. Construction of \mathcal{D} .

983
984 **R1 Seeds.** $\mathcal{D}^0 := \{0, e_1, \dots, e_{d_{\text{sty}}}\}$.

985
986 **R2 Two-of-three closure.** Let $k, k_1 \in \mathbb{Z}^{d_{\text{sty}}}$ satisfy $|k^{(i)}|, |(k+k_1)^{(i)}| \leq M_i$ for all i . Suppose
987 there exist latent points

$$988 z_A, z_B := z_A + k \odot q_0, z_C := z_A + (k + k_1) \odot q_0,$$

989 in \mathcal{Z} such that at least two of $\{z_A, z_B, z_C\}$ belong to the training set $\mathcal{Z}_{\text{train}}$, then, whenever
990 two vectors among $\{k, k_1, k+k_1\}$ already lie in the current set \mathcal{D}^n , insert the third to \mathcal{D}^{n+1}
991

992 Iterating **R2** yields an ascending chain $\mathcal{D}^0 \subset \mathcal{D}^1 \subset \dots$; abuse $\mathcal{D} := \bigcup_{n \geq 0} \mathcal{D}^n$.

995 **B. Training-time equivariance** For each augmentation $a_i(t_i)$, there is a transformation network
996 $g_i : \mathcal{T}_{a_i} \rightarrow [-\pi, \pi]$. Assume the training satisfies:

997
998 **E1 Cluster-wise perfect equivariance.** All the data are assigned to the corresponding cluster
999 determined by \mathcal{A} . For any index set $S \subseteq \{1, \dots, d_{\text{sty}}\}$, parameters $(t_i)_{i \in S}$, and all $z \in$
1000 $\mathcal{Z}_{\text{train}}$, calling only $\{a_i(t_i)\}_{i \in S}$ yields

$$1001 \tilde{f}(a_S(t, \varphi(z))) = \mathcal{R}\left(\sum_{i \in S} g_i(t_i) e_i\right) \tilde{f}(\varphi(z)),$$

1002
1003
1004 **E2 Non-zero canonical angle.** $\exists t_+^{(i)} \in \mathcal{T}_{a_i}$ such that $q_{a_i}(t_+^{(i)}) = q_0^{(i)}$ and $g_i(t_+^{(i)}) = \theta_0^{(i)} \neq 0$.

1005
1006 **E3 Slice non-degeneracy.** $\|\tilde{f}_{\text{sty}}^{(i)}(\varphi(z))\|_2 = r > 0$ for all $z \in \mathcal{Z}_{\text{train}}$ and all $i = 1, \dots, d_{\text{sty}}$.
1007

1008 **C. Conclusion.** (i) \mathcal{D} is well-defined and is the unique minimal subset satisfying **Seeds + Two-of-**
1009 **three closure.** (ii) For any $k \in \mathcal{D}$ one may form a parameter tuple $t(k) = (t_1, \dots, t_{d_{\text{sty}}})$ satisfying

$$1010 q_{a_i}(t_i) = \begin{cases} k^{(i)} q_0^{(i)}, & k^{(i)} \neq 0, \\ 0 & \text{(either by omitting } a_i \text{ or } t_i \text{ with } q_{a_i}(t_i) = 0), \quad k^{(i)} = 0, \end{cases}$$

1011
1012 and with this choice, write $g(t(k)) = (g_1(t_1), \dots, g_{\text{sty}}(t_{\text{sty}}))^T$

$$1013 g(t(k)) = k \odot \theta_0, \quad \theta_0 := (\theta_0^{(1)}, \dots, \theta_0^{(d_{\text{sty}})})^T.$$

1014
1015
1016
1017 **Intuition:** On the discrete set of style displacements actually "met" during training, the learned
1018 rotation angles are forced to be exactly linear in the semantic displacement. The linearity can be
1019 proved by propagating additivity from single-step moves to all reachable displacements.
1020

1021 *Proof.* (i) Rule **R2** always adds precisely the missing element of a triple $\{k, k_1, k+k_1\}$. The out-
1022 come after any finite sequence of rule applications therefore depends solely on the set of admissible
1023 triples, not on their order. Minimality follows by construction.

1024 (ii) We prove the linearity by induction over the depth at which k enters \mathcal{D} .

1025 Base.

- *Zero vector.* Take the empty call of augmentations (no a_i invoked). Perfect equivariance then supplies $\tilde{f}(\varphi(z)) = \mathcal{R}(g(t(0)))\tilde{f}(\varphi(z))$, since each $g_i(t_i)$ is in $[-\pi, \pi]$, we have $g(t(0)) = 0$.
- *Unit vector e_i* Activate *only* the i -th augmentation with its canonical parameter $t_+^{(i)}$ and leave every other augmentation uncalled: $a_{\{i\}}(t) = \{a_i(t_+^{(i)})\}$. Because $q_{a_i}(t_+^{(i)}) = q_0^{(i)}$ while all other displacements are zero, the latent displacement equals $e_i \odot q_0$. Perfect equivariance therefore yields

$$\tilde{f}(a_i(t_+^{(i)}), \varphi(z)) = \mathcal{R}_i(g_i(t_+^{(i)})) \tilde{f}(\varphi(z)),$$

and **E2** gives $g_i(t_+^{(i)}) = \theta_0^{(i)} \neq 0$. Consequently

$$g(t(e_i)) = (0, \dots, \theta_0^{(i)}, \dots, 0) = e_i \odot \theta_0.$$

Hence $g(t(k)) = k \odot \theta_0$ holds for every $k \in \mathcal{D}^0$.

Induction step. Assume $g(t(\ell)) = \ell \odot \theta_0$ for $\forall \ell \in \mathcal{D}^n$. Let **R2** add the previously missing vector of the triple $\{k, k_1, k + k_1\}$. Because the rule and the triple are permutation-symmetric, we may *without loss of generality* suppose $k, k_1 \in \mathcal{D}^n$, $k + k_1 \notin \mathcal{D}^n$. Likewise, the same symmetry lets us arrange the latent points so that $z_A, z_B \in \mathcal{Z}_{\text{train}}$ and $z_C \in \mathcal{Z}$. Augmentations are *only executed at training points* and each displacement of triplets is within the range of a single call of a (composite) augmentation, so all paths below are legal.

$$\begin{array}{lll} \text{Path I:} & z_A \xrightarrow{k \odot q_0} z_B & \tilde{f}(\varphi(z_B)) = \mathcal{R}(g(t(k))) \tilde{f}(\varphi(z_A)), \\ \text{Path II:} & z_B \xrightarrow{k_1 \odot q_0} z_C & \tilde{f}(\varphi(z_C)) = \mathcal{R}(g(t(k_1))) \tilde{f}(\varphi(z_B)), \\ \text{Path III:} & z_A \xrightarrow{(k+k_1) \odot q_0} z_C & \tilde{f}(\varphi(z_C)) = \mathcal{R}(g(t(k+k_1))) \tilde{f}(\varphi(z_A)), \end{array}$$

Composing Paths I and II and comparing with Path III, canceling the non-degenerated $\tilde{f}(\cdot)$, we should have the following matrix identity

$$\mathcal{R}(g(t(k+k_1))) = \mathcal{R}(g(t(k_1)))\mathcal{R}(g(t(k))).$$

Because $\theta \mapsto R_\theta$ is injective on $[-\pi, \pi]$ and the rotation blocks act on disjoint 2-planes (or add within the same plane), we deduce $g(t(k+k_1)) = g(t(k)) + g(t(k_1))$. Applying the induction hypothesis yields $g(t(k+k_1)) = (k+k_1) \odot \theta_0$. Therefore $g(t(\ell)) = \ell \odot \theta_0$ holds for every $\ell \in \mathcal{D}^{(n+1)}$, closing the induction. \square

Using the above lemma, we are ready to prove that, on the discrete grid of latent displacements seen in training, any two features in the same cluster differ only by a block-diagonal rotations whose angle vector is linear in their latent difference, as formalized in the following proposition.

Proposition 6a (Weak cluster-wise disentanglement on the discrete grid). *Assume*

A1 Discrete linearity of g (Lemma 2): for all $k \in \mathcal{D}$ there exist parameters $t(k)$ with $q_{a_i}(t_i) = k^{(i)} q_0^{(i)}$ such that $g(t(k)) = k \odot \theta_0$.

A2 Chain connectivity For any $z_A, z_B \in \mathcal{Z}_{\text{train}}$ with $z_A \sim_A z_B$, there exists a finite chain

$$z_A = z_0, z_1, \dots, z_n = z_B \quad (\text{all in } \mathcal{Z}_{\text{train}})$$

such that for each edge (z_j, z_{j+1}) there exist augmentations $u_j^+(t_j^+)$ and $u_j^-(t_j^-)$ with

$$\varphi^{-1}(u_j^+(t_j^+, \varphi(z_j))) = \varphi^{-1}(u_j^-(t_j^-, \varphi(z_{j+1}))) =: z_j^* \in \mathcal{Z},$$

and associated style displacements $k_j^+ \odot q^0$ and $k_j^- \odot q^0$, respectively, with $k_j^+, k_j^- \in \mathcal{D}$.

Then every cluster admits **weak disentanglement** on the discrete grid: for any two training points x_A, x_B in the same cluster,

$$\tilde{f}(\varphi(z_B)) = \mathcal{R}(k \odot \theta_0) \tilde{f}(\varphi(z_A)), \quad k \odot q_0 = z_B - z_A$$

1080 *Proof.* Fix an edge (z_j, z_{j+1}) of the chain in **A2**. By perfect equivariance and **A1** (which provides
1081 $g(t(k)) = k \odot \theta_0^{(i)}$ on \mathcal{D}), we obtain

$$1083 \quad \tilde{f}(\varphi(z_j^*)) = \mathcal{R}(k_j^+ \odot \theta_0) \tilde{f}(\varphi(z_j)), \quad \tilde{f}(\varphi(z_j^*)) = \mathcal{R}(k_j^- \odot \theta_0) \tilde{f}(\varphi(z_{j+1})).$$

1084 Equating and right-multiplying by the inverse rotation yields

$$1086 \quad \tilde{f}(\varphi(z_{j+1})) = \mathcal{R}((k_j^+ - k_j^-) \odot \theta_0) \tilde{f}(\varphi(z_j)).$$

1088 Composing these relations along $j = 0, \dots, n-1$, and using commutativity of block-diagonal
1089 rotations, gives

$$1091 \quad \tilde{f}(\varphi(z_B)) = \prod_{j=0}^{n-1} \mathcal{R}((k_j^+ - k_j^-) \odot \theta_0) \tilde{f}(\varphi(z_A)) = \mathcal{R}\left(\sum_{j=0}^{n-1} (k_j^+ - k_j^-) \odot \theta_0\right) \tilde{f}(\varphi(z_A)).$$

1094 By the construction of edges, $z_{j+1} - z_j = (k_j^+ - k_j^-) \odot q_0$. Summing over j gives

$$1096 \quad z_B - z_A = \sum_{j=0}^{n-1} (k_j^+ - k_j^-) \odot q_0 = k \odot q_0.$$

1099 Hence,

$$1100 \quad \tilde{f}(\varphi(z_B)) = \mathcal{R}(k \odot \theta_0) \tilde{f}(\varphi(z_A)).$$

1101 \square

1103 **Next, we will prove that the linearity succeed with continuous augmentation.**

1104 **Definition 5.** A realized displacement $s \in [-M_i, M_i]$ means there are $z_0, z_1 \in \mathcal{Z}_{\text{train}}$ with $z_1 - z_0 =$
1105 se_i and almost all intermediate states obtained by applying $a_i(t_i)$ to these two points remain inside
1106 \mathcal{Z} .

1107 **Lemma 3.** Let $\mathcal{A} = \{a_i(t_i)\}_{i=1}^{d_{\text{sty}}}$ be d_{sty} parametric augmentations, each acting solely on, without
1108 loss of generality, the i -th style coordinate and assigning to the i -th canonical rotation plane. For
1109 every latent $z \in \mathcal{Z}_{\text{train}}$ and parameter $t_i \in \mathcal{T}_{a_i}$ define

$$1111 \quad \Delta_{a_i}(t_i, z) := q_{a_i}(t_i) e_i, \text{ where } q_{a_i}(t_i) \in [-M_i, M_i]$$

1112 where e_i is the i -th standard basis vector of \mathbb{R}^m and $[-M_i, M_i] \subset \mathbb{R}$ is coordinate bound.

1114 Assume the training-time equivariance conditions **B** (items **E1** and **E3**) in Lemma 2, and each $g_i :$
1115 $\mathcal{T}_{a_i} \rightarrow [-\pi, \pi]$ and q_{a_i} are continuous. Write $\hat{g}_i(d) := g_i(q_{a_i}^{-1}(d))$, $d \in [-M_i, M_i]$, and assume
1116 $\hat{g}_i(d) \neq 0$ for any non-zero displacement. We have

$$1117 \quad 1. \hat{g}_i(0) = 0;$$

1118 2. For every realized displacement s , the one-step identity holds on admissible pairs:

$$1120 \quad \hat{g}_i(d) - \hat{g}_i(d - s) = \theta_s \quad \text{whenever } d, d - s \in [-M_i, M_i], \quad \theta_s := \hat{g}_i(s) \neq 0. \quad (14)$$

1122 *Proof.* For any $z \in \mathcal{Z}_{\text{train}}$, Condition **E1** gives $\tilde{f}(\varphi(z)) = \mathcal{R}_i(\hat{g}_i(0)) \tilde{f}(\varphi(z))$, since each $g_i(t_i)$ is
1123 in $[-\pi, \pi]$, we have $\hat{g}_i(0) = 0$. Thus prove (1).

1124 For any realized displacement s , fix an arbitrary $d \in [-M_i, M_i]$ with $d - s \in [-M_i, M_i]$, we have:

$$1127 \quad \begin{aligned} \text{Path I:} & \quad z_0 \xrightarrow{se_i} z_1 & \tilde{f}(\varphi(z_1)) &= \mathcal{R}_i(\hat{g}_i(r)) \tilde{f}(\varphi(z_0)), \\ \text{Path II:} & \quad z_1 \xrightarrow{(d-s)e_i} z_2 & \tilde{f}(\varphi(z_2)) &= \mathcal{R}_i(\hat{g}_i(d-s)) \tilde{f}(\varphi(z_1)), \\ \text{Path III:} & \quad z_0 \xrightarrow{de_i} z_2 & \tilde{f}(\varphi(z_2)) &= \mathcal{R}_i(\hat{g}_i(d)) \tilde{f}(\varphi(z_0)). \end{aligned}$$

1132 Composing Paths I and II and comparing with Path III gives

$$1133 \quad \mathcal{R}_i(\hat{g}_i(d)) \tilde{f}(\varphi(z_0)) = \mathcal{R}_i(\hat{g}_i(d-s)) \mathcal{R}_i(\hat{g}_i(r)) \tilde{f}(\varphi(z_0))$$

By the condition **E3**, $\tilde{f}(\varphi(z_0))$ is non-degenerated, we have

$$\mathcal{R}_i(\widehat{g}_i(d)) = \mathcal{R}_i(\widehat{g}_i(d-s))\mathcal{R}_i(\widehat{g}_i(r)) = \mathcal{R}_i(\widehat{g}_i(d-s) + \widehat{g}_i(r)).$$

By injectivity of the rotation map, this yields $\widehat{g}_i(d) - \widehat{g}_i(d-s) = \theta_s$ with $\theta_s := \widehat{g}_i(s) \neq 0$. \square

Lemma 4. *Under the assumption of Lemma 3. Let $S \subset [-M_i, M_i]$ be a nonempty family of realized steps that is closed under admissible common multiples: for any $r, s \in S$, there exist $m, n \in \mathbb{Z}$ such that $L := mr = ns \in [-M_i, M_i]$. Then the quantity*

$$c_S := \frac{\theta_s}{s} \neq 0 \quad (s \in S)$$

is well-defined (independent of $s \in S$), and

$$\widehat{g}_i(d) = c_S d + H_S(d), \quad H_S(d-s) = H_S(d) \quad (\forall s \in S),$$

on every subinterval where the pairs in Equation equation 14 are admissible. In particular, H_S is s -periodic for each $s \in S$.

Proof. Apply the three paths in Lemma 3 at the sequence $d_j := jr$ ($j = 1, \dots, m$), all of which lie in $[-M_i, M_i]$:

$$\widehat{g}_i(d_j) - \widehat{g}_i(d_{j-1}) = \theta_r \quad \Rightarrow \quad \widehat{g}_i(L) - \widehat{g}_i(0) = \sum_{j=1}^m \theta_r = m \theta_r.$$

Similarly, using step s at $e_\ell := \ell s$ ($\ell = 1, \dots, n$),

$$\widehat{g}_i(L) - \widehat{g}_i(0) = n \theta_s.$$

Therefore

$$\frac{\theta_r}{r} = \frac{\theta_s}{s} =: c_S.$$

So the slope is independent of the realized step. Define the common remainder $H_S(d) := \widehat{g}_i(d) - c_S d$.

$$H_S(d) - H_S(d-s) = (\widehat{g}_i(d) - \widehat{g}_i(d-s)) - \frac{\widehat{g}_i(s)}{s}(d - (d-s)) = \widehat{g}_i(s) - \widehat{g}_i(s) = 0.$$

Therefore

$$\widehat{g}_i(d) = c_S d + H_S(d), \quad H_S(d-s) = H_S(d) \quad (\forall s \in S). \quad (15)$$

\square

Lemma 5. *Under the conditions of Lemma 3 and Lemma 4. Let $S \subset [-M_i, M_i]$ be a nonempty common multiplier-closed family of realized displacements, so that $c_S = \theta_s/s$ is well defined for all $s \in S$ and $H_S(d) := \widehat{g}_i(d) - c_S d$ is s -periodic on admissible pairs for every $s \in S$. Suppose moreover that $\widehat{g}_i : [-M_i, M_i] \rightarrow \mathbb{R}$ be L -Lipschitz:*

$$|\widehat{g}_i(x) - \widehat{g}_i(y)| \leq L|x-y| \quad \forall x, y \in [-M_i, M_i].$$

(a) **Uniform small remainder from the smallest step in a family.** Assume S contains the smallest magnitude step

$$s_* := \arg \min_{s \in S} |s|.$$

Then, using only the periodicity provided by Lemma 4 for the family S ,

$$|H_S(d)| \leq 2L|s_*|.$$

(b) **Near-loop slope matching across two families.** Let $S_1, S_2 \subset [-M_i, M_i]$ be two common multiplier-closed families with slopes c_{S_1}, c_{S_2} (from Lemma 4). Suppose there exist $r \in S_1$, $s \in S_2$ and integers $m, n \neq 0$ such that, for some basepoint $d \in [-M_i, M_i]$, the chains

$$d \rightarrow d + mr, \quad d \rightarrow d + ns$$

are admissible and the endpoints are δ -close: $|mr - ns| \leq \delta$. Then

$$|c_{S_1} - c_{S_2}| \leq \frac{2L\delta}{\max\{m|r|, n|s|\}}$$

and, for any $x \in [-M_i, M_i]$,

$$|H_{S_1}(x) - H_{S_2}(x)| \leq M_i |c_{S_1} - c_{S_2}| \leq \frac{2LM_i\delta}{\max\{m|r|, n|s|\}}.$$

Moreover, each family remainder is $2L\delta$ -almost periodic across the other family's step:

$$|H_{S_1}(d + ns) - H_{S_1}(d)| \leq 2L\delta, \quad |H_{S_2}(d + mr) - H_{S_2}(d)| \leq 2L\delta,$$

Proof. By Lemma 4 the quotient $c_S = \theta_s/s$ is fixed for any $s \in S$, and $H_S(d) = \widehat{g}(d) - c_S d$ is s -periodic for each $s \in S$ on admissible pairs. Moreover, since \widehat{g} is L -Lipschitz and $\theta_s = \widehat{g}_i(s) - 0 = \widehat{g}(s) - \widehat{g}(0)$,

$$|c_S| = \left| \frac{\theta_s}{s} \right| = \frac{|\widehat{g}(s) - \widehat{g}(0)|}{|s|} \leq L \quad \text{for any } s \in S. \quad (*)$$

(a) **Uniform bound from the smallest step.** Fix $d \in [-M_i, M_i]$. By s_* -periodicity (iterated along the segment joining d to $d - ks_* \in [-M_i, M_i]$), choose $k \in \mathcal{Z}$ such that

$$d^\circ := d - ks_* \in [-|s_*|, |s_*|] \cap [-M_i, M_i] \quad \text{and} \quad H_S(d) = H_S(d^\circ).$$

Then, using $H_S(0) = \widehat{g}(0)$ and equation *,

$$\begin{aligned} |H_S(d)| &= |H_S(d^\circ)| = |\widehat{g}(d^\circ) - c_S d^\circ| \leq |\widehat{g}(d^\circ) - \widehat{g}(0)| + |c_S| |d^\circ| \\ &\leq L |d^\circ| + L |d^\circ| \leq 2L |s_*|. \end{aligned}$$

This proves the claimed uniform bound.

(b) **Near-loop slope matching and cross-almost-periodicity.** By admissibility and additivity along the two chains,

$$\widehat{g}(d + mr) = \widehat{g}(d) + m\theta_r = \widehat{g}(d) + m c_{S_1} r, \quad \widehat{g}(d + ns) = \widehat{g}(d) + n\theta_s = \widehat{g}(d) + n c_{S_2} s.$$

Subtract and use the L -Lipschitz property:

$$|m c_{S_1} r - n c_{S_2} s| = |\widehat{g}(d + mr) - \widehat{g}(d + ns)| \leq L |mr - ns| \leq L\delta.$$

By the triangle inequality,

$$\begin{aligned} |m c_{S_1} r - n c_{S_2} s| &= |mr(c_{S_1} - c_{S_2}) + c_{S_2}(mr - ns)| \\ &\geq |mr| |c_{S_1} - c_{S_2}| - |c_{S_2}| |mr - ns| \end{aligned}$$

Then,

$$\begin{aligned} |c_{S_1} - c_{S_2}| &\leq \frac{|m c_{S_1} r - n c_{S_2} s| + |c_{S_2}| |mr - ns|}{m|r|} \\ &\leq \frac{L\delta + |c_{S_2}| |mr - ns|}{m|r|} \\ &\leq \frac{L\delta + L\delta}{m|r|} \quad \text{Using } (*) \\ &= \frac{2L\delta}{m|r|} \end{aligned}$$

Interchanging the roles of (r, m) and (s, n) gives

$$|c_{S_1} - c_{S_2}| \leq \frac{2L\delta}{n|s|}.$$

Combining the two bounds yields the displayed estimate for $|c_{S_1} - c_{S_2}|$,

$$|c_{S_1} - c_{S_2}| \leq \frac{2L\delta}{\max\{m|r|, n|s|\}}$$

For the remainders,

$$|H_{S_1}(x) - H_{S_2}(x)| = |(c_{S_2} - c_{S_1})x| \leq |x| |c_{S_1} - c_{S_2}| \leq M_i |c_{S_1} - c_{S_2}|,$$

which gives the second display.

Finally, for the cross-almost-periodicity, note that

$$\begin{aligned} |H_{S_1}(d + ns) - H_{S_1}(d)| &= |\widehat{g}(d + ns) - \widehat{g}(d) - c_{S_1} ns| \\ &= |n c_{S_2} s - c_{S_1} ns| = n|s| |c_{S_2} - c_{S_1}|. \end{aligned}$$

Since H_{S_1} is r -periodic and $H_{S_1}(d + mr) = H_{S_1}(d)$

$$|H_{S_1}(d + ns) - H_{S_1}(d)| = |H_{S_1}(d + ns) - H_{S_1}(d + mr)|.$$

Using the previously derived bound with the triangle trick, we have

$$\begin{aligned} |H_{S_1}(d + ns) - H_{S_1}(d)| &= |\widehat{g}(d + ns) - \widehat{g}(d + mr) - c_{S_1}(ns - mr)| \\ &\leq |\widehat{g}(d + ns) - \widehat{g}(d + mr)| + |c_{S_1}| |ns - mr| \\ &\leq L\delta + L\delta = 2L\delta. \end{aligned}$$

The estimate for H_{S_2} is identical by symmetry. This completes the proof. \square

Lemma 3 simply shows that along a fixed style coordinate, each realized step s must always induce the same angular increment θ_s , independent of where you start. Lemma 4 then says that, within any family S of steps that share common multiples, these increments are all compatible: they define a single common slope c_S , and \widehat{g}_i decomposes into a linear part $c_S d$ plus an S -periodic remainder H_S . Finally, Lemma 5 uses Lipschitz regularity to show that this periodic remainder cannot oscillate much: small steps force H_S to be uniformly small, and “almost closed loops” built from two step families force their slopes and remainders to almost coincide. In combination, these lemmas say that on a dense network of realized displacements, the transformation network g_i is essentially linear in latent $z^{(i)}$ with only a very small oscillatory part.

Corollary 1 (Linearity from either an irrational pair or sub-step refinement). *Fix a coordinate i . Assume Lemma 3 (one-step identity/additivity on admissible chains) and Lemma 4 (common slope on a common multiplier-closed family of realized displacements), so that for any realized step $r > 0$ we have*

$$c_i = \frac{\theta_r}{r} \neq 0 \quad (\text{independent of } r), \quad H_i(d) := \widehat{g}_i(d) - c_i d$$

and H_i is r -periodic on admissible pairs. If, in addition, either

(i) there exist realized displacements $r_i \in (0, M_i]$ and $\beta_i r_i \in (0, M_i]$ with $\beta_i \in (0, 1) \setminus \mathbb{Q}$, or

(ii) there exists a sequence of realized displacements $r_n = \frac{r_0}{n}$ with $r_n \rightarrow 0$,

then $H_i \equiv 0$ on $[-M_i, M_i]$ and hence

$$\widehat{g}_i(d) = c_i d \quad \text{for all } d \in [-M_i, M_i], \quad g_i(t_i) = c_i q_{a_i}(t_i) \quad \text{for all } t_i \in \mathcal{T}_{a_i}.$$

1296 *Proof.* We give separate arguments.
1297

1298 *Case (i): incommensurate pair r_i and $\beta_i r_i$.* Fix $r_i \in (0, M_i]$ realized and $\beta_i \in (0, 1) \setminus \mathbb{Q}$ such that
1299 $\beta_i r_i$ is realized. From Equation 15 with $s = r_i$,

$$1300 H_i(d) = H_i(d - r_i) \quad \text{whenever } d, d - r_i \in [-M_i, M_i].$$

1302 Next, derive the precise $\beta_i r_i$ -increment identity for H_i . For any l with $l, l - \beta_i r_i \in [-M_i, M_i]$,

$$\begin{aligned} 1303 H_i(l - \beta_i r_i) &= \widehat{g}_i(l - \beta_i r_i) - \frac{\theta_0^{(i)}}{r_i} (l - \beta_i r_i) \\ 1304 &= (\widehat{g}_i(l) - \widehat{g}_i(\beta_i r_i)) - \frac{\theta_0^{(i)}}{r_i} l + \theta_0^{(i)} \beta_i \\ 1305 &= \left(\widehat{g}_i(l) - \frac{\theta_0^{(i)}}{r_i} l \right) - \left(\widehat{g}_i(\beta_i r_i) - \theta_0^{(i)} \beta_i \right) \\ 1306 &= H_i(l) - C, \end{aligned} \tag{16}$$

1312 where we set

$$1313 C := \widehat{g}_i(\beta_i r_i) - \theta_0^{(i)} \beta_i \quad (\text{a constant independent of } l).$$

1315 Since $r_i \leq M_i$, it suffices to analyze H_i on $I := (0, r_i] \subset [-M_i, M_i]$ and extend by r_i -periodicity.
1316 Define the sequence $\{d_m\}_{m \geq 1} \subset (0, r_i]$ by

$$1317 d_1 := r_i - \beta_i r_i, \quad d_{m+1} := \begin{cases} d_m - \beta_i r_i, & d_m \geq \beta_i r_i, \\ d_m - \beta_i r_i + r_i, & d_m < \beta_i r_i. \end{cases}$$

1320 Claim 1. For all $m \geq 1$,

$$1321 H_i(d_m) = -mC, \quad d_m \in (0, r_i]. \tag{17}$$

1322 *Proof of Claim 1.* Base case $m = 1$: taking $l = r_i$ in equation 16 (both r_i and $r_i - \beta_i r_i$ lie in I),

$$1323 H_i(d_1) = H_i(r_i - \beta_i r_i) = H_i(r_i) - C = 0 - C = -C.$$

1325 Induction step $m \rightarrow m + 1$: If $d_m \geq \beta_i r_i$, then $d_{m+1} = d_m - \beta_i r_i \in (0, r_i]$ and equation 16 with
1326 $l = d_m$ gives $H_i(d_{m+1}) = H_i(d_m) - C$. If instead $d_m < \beta_i r_i$, then $d_{m+1} = d_m - \beta_i r_i + r_i \in (0, r_i]$
1327 and equation 16 yields $H_i(d_m) - H_i(d_m - \beta_i r_i) = C$; by r_i -periodicity, $H_i(d_m - \beta_i r_i) = H_i(d_m -$
1328 $\beta_i r_i + r_i) = H_i(d_{m+1})$, hence again $H_i(d_{m+1}) = H_i(d_m) - C$. This proves equation 17.
1329

1330 Claim 2. The sequence $\{d_m\}$ is dense in $(0, r_i]$.

1331 *Proof of Claim 2.* By construction,

$$1332 d_{m+1} \equiv d_m - \beta_i r_i \pmod{r_i}, \quad \text{so } d_m \equiv r_i - m\beta_i r_i \pmod{r_i}.$$

1333 To pass from a congruence class to its representative in $(0, r_i]$, use the identity

$$1334 x \bmod r_i = x - r_i \lfloor \frac{x}{r_i} \rfloor = r_i \left\{ \frac{x}{r_i} \right\} \quad \text{for any } x \in \mathbb{R},$$

1335 where $\lfloor x \rfloor$ is the greatest integer less than or equal to x and $\{x\} = x - \lfloor x \rfloor$ is the fractional part
1336 in $[0, 1)$. Applying this to $x = r_i - m\beta_i r_i = r_i(1 - m\beta_i)$ and using the elementary relation
1337 $\{1 - y\} = 1 - \{y\}$ for $y \notin \mathbb{Z}$ (which holds here because $\beta_i \notin \mathbb{Q}$ implies $m\beta_i \notin \mathbb{Z}$ for all $m \geq 1$),
1338 we obtain

$$1339 d_m \equiv (r_i - m\beta_i r_i) \bmod r_i = r_i \left\{ \frac{r_i - m\beta_i r_i}{r_i} \right\} = r_i \{1 - m\beta_i\} = r_i (1 - \{m\beta_i\}). \tag{18}$$

1340 Since $\beta_i \notin \mathbb{Q}$, the set $\{\{m\beta_i\} : m \in \mathbb{N}\}$ is dense in $[0, 1)$ (Kronecker's theorem), hence $\{d_m\}$ is
1341 dense in $(0, r_i]$ by the continuity of $x \mapsto r_i(1 - x)$ on $[0, 1)$. In particular, there is a subsequence
1342 $d_{m_\ell} \rightarrow 0^+$.

1343 From equation 17, $H_i(d_{m_\ell}) = -m_\ell C$. Since $d_{m_\ell} \rightarrow 0^+$ and H_i is continuous with $H_i(0) = 0$,
1344 we have $H_i(d_{m_\ell}) \rightarrow 0$, hence $-m_\ell C \rightarrow 0$. As $m_\ell \rightarrow \infty$, the only possibility is $C = 0$. With
1345

1350 $C = 0$, equation 16 gives $H_i(l - \beta_i r_i) = H_i(l)$ on admissible pairs, so H_i is both r_i - and $\beta_i r_i$ -
 1351 periodic. From equation 17 we also have $H_i(d_m) = 0$ for all m , so H_i vanishes on the dense set
 1352 $\{d_m\} \subset (0, r_i]$; by continuity, $H_i \equiv 0$ on $(0, r_i]$, and then by r_i -periodicity $H_i \equiv 0$ on $[-M_i, M_i]$.
 1353

1354 *Case (ii): periods accumulating at 0.* By hypothesis, for each N large enough the sub-step $r_N :=$
 1355 r_0/N is realized and H_i is r_N -periodic on admissible pairs by Lemma 4. Since $r_N \rightarrow 0$, we prove
 1356 H_i is constant on $[-M_i, M_i]$ by a uniform-continuity argument.
 1357

1358 Because H_i is continuous on the compact interval $[-M_i, M_i]$, it is uniformly continuous: there exists
 1359 a modulus $\omega(\cdot)$ with $\omega(\delta) \rightarrow 0$ as $\delta \rightarrow 0$ and $|H_i(x) - H_i(y)| \leq \omega(|x - y|)$ for all $x, y \in [-M_i, M_i]$.
 1360 Fix arbitrary $x, y \in [-M_i, M_i]$ and $\varepsilon > 0$. Choose N so large that $r_N < \varepsilon$ and $\omega(r_N) < \varepsilon$, and
 1361 such that the r_N -periodicity applies to every pair along the short segment we construct next. Pick
 1362 $m \in \mathbb{Z}$ with $|x - (y + mr_N)| \leq r_N$ and with the points $y, y + r_N, \dots, y + mr_N$ all in $[-M_i, M_i]$
 1363 (possible since r_N is small). Then $H_i(y + mr_N) = H_i(y)$ by r_N -periodicity, and

$$1364 |H_i(x) - H_i(y)| = |H_i(x) - H_i(y + mr_N)| \leq \omega(|x - (y + mr_N)|) \leq \omega(r_N) < \varepsilon.$$

1365 Since $\varepsilon > 0$ and x, y are arbitrary, H_i is constant on $[-M_i, M_i]$; evaluating at 0 gives $H_i \equiv 0$.
 1366

1367 In both cases, $H_i \equiv 0$ on $[-M_i, M_i]$, hence $\widehat{g}_i(d) = c_i d$ for all $d \in [-M_i, M_i]$, and substituting $d =$
 1368 $q_{a_i}(t_i)$ yields $g_i(t_i) = c_i q_{a_i}(t_i)$ for all $t_i \in \mathcal{T}_{a_i}$. Apply the same argument to all $i \in [1, \dots, d_{\text{sty}}]$,
 1369 we can get the linearity of $g_i(t_i)$ for all i . \square
 1370

1371 **Intuition:** The previous lemmas (Lemma 3, 4, 5) tell us that, along a single style coordinate, the
 1372 learned rotation angle can be written as a linear function plus a periodic remainder that repeats with
 1373 every realized step. Corollary 1 uses an extra richness assumption on the realized steps (either an
 1374 irrational pair of step sizes, or steps that can be refined to arbitrarily small increments) to show that
 1375 such a continuous, multi-periodic remainder must in fact be constant. Since the remainder is zero
 1376 at zero displacement, it must be zero everywhere, so the angle becomes exactly linear in the latent
 1377 displacement.

1378 **Remark 2** (Empirical remainder and quantitative guarantees). *In our experiments the measured*
 1379 *remainder $H_i(d) = \widehat{g}_i(d) - c_i d$ is numerically ≈ 0 , indicating that the learned g_i are effectively*
 1380 *linear. This aligns not only with the qualitative routes in Corollary 1, but also with the quantitative*
 1381 *bounds of Proposition 5:*

- 1382 (a) (Uniform small remainder from the smallest realized step) *If the realized family contains a*
 1383 *smallest magnitude step $s_* > 0$, then*

$$1384 \sup_{d \in [-M_i, M_i]} |H_i(d)| \leq 2L |s_*|.$$

1385 Hence, as training realizes finer sub-steps ($|s_*| \rightarrow 0$), the remainder vanishes uniformly.
 1386

- 1387 (b) Near-loop self-consistency within the same augmentation. *For any realized step sizes r, s and*
 1388 *integers $m, n \geq 1$, if there is a basepoint $d \in [-M_i, M_i]$ such that the chains $d \rightarrow d + mr$ and*
 1389 *$d \rightarrow d + ns$ are admissible and their endpoints are δ -close ($|mr - ns| \leq \delta$), then*

$$1390 \left| \frac{\theta_r}{r} - \frac{\theta_s}{s} \right| \leq \frac{2L \delta}{\max\{m|r|, n|s|\}}, \quad |H_i(d + ns) - H_i(d + mr)| \leq 2L \delta.$$

1391 Thus, tighter near-loops (δ small, m, n large) force the inferred slopes to agree and make the
 1392 remainder nearly periodic across different realized steps.
 1393

1394 These quantitative effects explain why we observe $H_i \approx 0$ in practice under a single continuous
 1395 augmentation: sub-step refinement drives (a) to zero, and even when only approximate rational
 1396 relations are realized, (b) keeps the learned slopes and remainders tightly controlled.
 1397

1398 **Proposition 6b** (Weak cluster-wise disentanglement in the continuous case). *Assume*

1399 **B1 Linearity of g (Colloary 1).** *For $\forall i = [1, \dots, d_{\text{sty}}]$, $\forall t_i \in \mathcal{T}_{a_i}$, $g_i(t_i) = c_i q_{a_i}(t_i)$, write*
 1400 $g(t) := (g_1(t_1), \dots, g_{d_{\text{sty}}}(t_{d_{\text{sty}}}))^T$, $c_0 := (c_1, \dots, c_{d_{\text{sty}}})^T$.
 1401

B2 Chain-connectivity For any $z_A, z_B \in \mathcal{Z}_{\text{train}}$, with $z_A \sim_{\mathcal{A}} z_B$, there exists a finite chain

$$z_A = z_0, z_1, \dots, z_n = z_B \quad (\text{all in } \mathcal{Z}_{\text{train}})$$

such that for each edge (z_j, z_{j+1}) there are two augmentations $u_j^+(t_j^+), u_j^-(t_j^-)$ and a latent point $z_j^* \in \mathcal{Z}$ with

$$\varphi^{-1}(u_j^+(t_j^+, \varphi(z_j))) = \varphi^{-1}(u_j^-(t_j^-, \varphi(z_{j+1}))) =: z_j^*,$$

and realized displacements $m_j^+ := u_j^+(t_j^+)$, $m_j^- := u_j^-(t_j^-) \in \mathbb{R}^{d_{\text{sty}}}$ satisfying $m_j^{+, (i)}, m_j^{-, (i)} \in [-M_i, M_i]$ for all i .

Then every cluster admits **weak disentanglement**: for any two training points x_A, x_B in the same cluster,

$$\tilde{f}(\varphi(z_B)) = \mathcal{R}(m \odot c_0) \tilde{f}(\varphi(z_A)), \quad m = z_B - z_A$$

Proof. Since perfect equivariance is achieved, for a fixed edge (z_j, z_{j+1}) we have

$$\tilde{f}(\varphi(z_j^*)) = \mathcal{R}(g(u_j^+)) \tilde{f}(\varphi(z_j)) = \mathcal{R}(m_j^+ \odot c_0) \tilde{f}(\varphi(z_j)),$$

and also

$$\tilde{f}(\varphi(z_j^*)) = \mathcal{R}(g(u_j^-)) \tilde{f}(\varphi(z_{j+1})) = \mathcal{R}(m_j^- \odot c_0) \tilde{f}(\varphi(z_{j+1})).$$

Canceling the common term yields

$$\tilde{f}(\varphi(z_{j+1})) = \mathcal{R}(m_j^+ - m_j^-) \odot \theta_0 \tilde{f}(\varphi(z_j)).$$

Compose over $j = 0, \dots, n-1$ and use commutativity of block rotations:

$$\tilde{f}(\varphi(z_B)) = \mathcal{R}\left(\sum_{j=0}^{n-1} (m_j^+ - m_j^-) \odot c_0\right) \tilde{f}(\varphi(z_A)).$$

By construction of each edge, $z_{j+1} - z_j = m_j^+ - m_j^-$, hence $m = z_B - z_A = \sum_{j=0}^{n-1} (m_j^+ - m_j^-)$, which implies that:

$$\tilde{f}(\varphi(z_B)) = \mathcal{R}(m \odot c_0) \tilde{f}(\varphi(z_A))$$

□

Intuition: This proof is similar to Proposition 6a.

C.4 PROOF OF PROPOSITION 3

Definition 6 (Canonical nonnegative feature angle on a style plane). *Fix the counterclockwise orientation of each style plane. For the i -th plane, write $\tilde{f}_{\text{sty}}^{(i)}(x) = (u_i(x), v_i(x)) \in \mathbb{S}^1 \subset \mathbb{R}^2$. and two connected points $\varphi(z_A) = x_A \sim_{\mathcal{A}} x_B = \varphi(z_B)$ with $v := z_B - z_A$, denote the principal relative angle by*

$$\begin{aligned} \theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)) &:= \text{atan2}\left(\tilde{f}_{\text{sty}}^{(i)}(x_A) \times \tilde{f}_{\text{sty}}^{(i)}(x_B), \tilde{f}_{\text{sty}}^{(i)}(x_A) \cdot \tilde{f}_{\text{sty}}^{(i)}(x_B)\right) \\ &= \text{atan2}\left(u_i(x_A) v_i(x_B) - v_i(x_A) u_i(x_B), u_i(x_A) u_i(x_B) + v_i(x_A) v_i(x_B)\right) \\ &\in (-\pi, \pi]. \end{aligned}$$

Define the canonical feature angle

$$\Theta_i^{\text{ccw}}(\tilde{f}(x_B); \tilde{f}(x_A)) := \begin{cases} \theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)), & \theta_i^{\text{pr}} \geq 0, \\ \theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)) + 2\pi, & \theta_i^{\text{pr}} < 0, \end{cases} \in [0, 2\pi),$$

and the directed feature angle (follow the sign of $v^{(i)}$) by

$$\Theta_i^{\text{dir}}(\tilde{f}(x_B); \tilde{f}(x_A); v^{(i)}) := \begin{cases} \Theta_i^{\text{ccw}}(\tilde{f}(x_B); \tilde{f}(x_A)) \in [0, 2\pi), & v^{(i)} \geq 0, \\ \Theta_i^{\text{ccw}}(\tilde{f}(x_B); \tilde{f}(x_A)) - 2\pi \in (-2\pi, 0], & v^{(i)} < 0. \end{cases}$$

Proposition 7 (Strong cluster-wise disentanglement). *Let $\mathcal{A} = \{a_i(t_i)\}_{i=1}^{d_{\text{sty}}}$ be d_{sty} parametric augmentations, each acting solely on the, without loss of generality, i -th style coordinate and assigning to i th canonical rotation plane. Under the conditions:*

S1 Training-time equivariance with linear head (Corollary 1 or Lemma 2). *There exist slopes $\theta_0^{(i)}/r_i \neq 0$ (for discrete case $r_i = 1$) such that any induced displacement v by augmentation output angles $g(t) = c_0 \odot v$. i.e., the i -th style plane rotates by $c_0^{(i)} v^{(i)}$. W.l.o.g, assume $c_0^{(i)} > 0$ for all i .*

S2 Connectivity (Proposition. 6b or Prop. 6a). *Any pair $\varphi(z_A) \sim_{\mathcal{A}} \varphi(z_B)$ is chain-connected with latent displacement $v := z_B - z_A$.*

Fix the (counterclockwise) orientation per the style plane and let $\Theta_i^{[0,2\pi]}(\cdot; \cdot)$ denote the canonical nonnegative feature angle (Def. 6). Then, for each style plane i ,

$$\Theta_i^{\text{dir}}(\tilde{f}(x_B); \tilde{f}(x_A); v^{(i)}) \equiv (c_0^{(i)} v^{(i)}) \bmod 2\pi \in [0, 2\pi). \quad (19)$$

If moreover the no-wrap budget holds, i.e., $|c_0^{(i)} v^{(i)}| < 2\pi$ for all $v^{(i)}$ on the data support, then no wrap occurs and the modulo in equation 19 disappears:

$$\Theta_i^{\text{dir}}(\tilde{f}(x_B); \tilde{f}(x_A); v^{(i)}) = c_0^{(i)} v^{(i)} \quad \left(\in [0, 2\pi) \text{ if } v^{(i)} \geq 0, \in (-2\pi, 0] \text{ if } v^{(i)} < 0 \right). \quad (20)$$

*Consequently, on each style plane i the directed feature angle depends only on $v^{(i)}$ and varies strictly monotonically with it, achieving **strong disentanglement**.*

Proof. Fix a style plane i and a connected pair

$$x_A = \varphi(z_A) \sim_{\mathcal{A}} x_B = \varphi(z_B), \quad v := z_B - z_A, \quad \theta := c_0^{(i)} v^{(i)}.$$

By **S1** (equivariance with a linear head) and **S2** (connectivity), the i -plane feature rotates by angle θ :

$$\tilde{f}_{\text{sty}}^{(i)}(x_B) = R(\theta) \tilde{f}_{\text{sty}}^{(i)}(x_A), \quad R(\theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}. \quad (21)$$

Write $a := \tilde{f}_{\text{sty}}^{(i)}(x_A) = (a_1, a_2)$ and $b := \tilde{f}_{\text{sty}}^{(i)}(x_B) = (b_1, b_2)$. From equation 21, either by direct multiplication or by writing $a = (\cos \alpha, \sin \alpha)$ and $b = (\cos(\alpha + \theta), \sin(\alpha + \theta))$, we obtain

$$d := a \cdot b = \cos \theta, \quad s := a \times b = \sin \theta. \quad (22)$$

By Definition 6, the principal relative angle is

$$\theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)) = \text{atan2}(s, d) = \text{atan2}(\sin \theta, \cos \theta) \in (-\pi, \pi].$$

By the defining property of atan2 ,

$$\theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)) \equiv \theta \pmod{2\pi} \quad (23)$$

By the same definition,

$$\Theta_i^{\text{ccw}}(\tilde{f}(x_B); \tilde{f}(x_A)) = \begin{cases} \theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)), & \theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)) \geq 0, \\ \theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)) + 2\pi, & \theta_i^{\text{pr}}(\tilde{f}(x_B); \tilde{f}(x_A)) < 0, \end{cases} \in [0, 2\pi).$$

Combining with equation 23 yields

$$\Theta_i^{\text{ccw}} \equiv \theta \pmod{2\pi} \quad (24)$$

By Definition 6,

$$\Theta_i^{\text{dir}}(\tilde{f}(x_B); \tilde{f}(x_A); v^{(i)}) = \begin{cases} \Theta_i^{\text{ccw}}(\tilde{f}(x_B); \tilde{f}(x_A)), & v^{(i)} \geq 0, \\ \Theta_i^{\text{ccw}}(\tilde{f}(x_B); \tilde{f}(x_A)) - 2\pi, & v^{(i)} < 0, \end{cases}$$

and equation 24 immediately gives the congruence

$$\Theta_i^{\text{dir}}(\tilde{f}(x_B); \tilde{f}(x_A); v^{(i)}) \equiv \theta = c_0^{(i)} v^{(i)} \pmod{2\pi}, \quad (25)$$

which is precisely equation 19.

Assume $|\theta| = |c_0^{(i)} v^{(i)}| < 2\pi$. Then:

- If $v^{(i)} \geq 0$ (hence $\theta \geq 0$), equation 24 gives $\Theta_i^{\text{ccw}} = \theta \in [0, 2\pi)$, hence $\Theta_i^{\text{dir}} = \Theta_i^{\text{ccw}} = \theta = c_0^{(i)} v^{(i)}$.
- If $v^{(i)} < 0$ (hence $\theta < 0$), equation 24 gives $\Theta_i^{\text{ccw}} = \theta + 2\pi \in (0, 2\pi)$, hence $\Theta_i^{\text{dir}} = \Theta_i^{\text{ccw}} - 2\pi = (\theta + 2\pi) - 2\pi = \theta = c_0^{(i)} v^{(i)} \in (-2\pi, 0]$.

This proves equation 20.

On the no-wrap domain, $\Theta_i^{\text{dir}} = c_0^{(i)} v^{(i)}$ with $c_0^{(i)} > 0$ (by our orientation convention), so $v^{(i)} \mapsto \Theta_i^{\text{dir}}$ is strictly increasing and does not depend on any $v^{(j)}$ with $j \neq i$. This is the stated strong disentanglement. If some $c_0^{(i)} < 0$, flip the orientation of the i -th style plane (swap its axes), which replaces $c_0^{(i)}$ by $|c_0^{(i)}|$ and preserves all statements. \square

Intuition: As the previous lemmas and propositions says, along each style coordinate, changing i -th latent $z^{(i)}$ by augmentation always rotates the i -th style plane by a fixed slope times that change. In this proposition we just read that angle back from the features themselves: the two feature points lie on the unit circle, so their relative angle is determined by a standard geometric construction (dot/cross products and atan2). If rotations never "wrap around" a full turn, there is no 2π ambiguity, so the measured feature angle is exactly this slope times augmentation induced displacement $v^{(i)}$.

Remark 3 (Angle budget). To enforce the budget $|c_0^{(i)} v^{(i)}| \leq 2\pi$ during training, constrain angles via Equation 13 in the main content and pick a safety factor η :

- **Cyclic Style Latents:** set $\eta = 1$ (full branch).
- **Non-cyclic Style Latents:** use a conservative $\eta < 1$ (e.g. $0.4 \sim 0.7$) to keep angles away from the endpoint and avoid wrap.
- **Approximate strength:** if $\max_{t \in T_{\text{data}}} \tilde{a}_a(t)$ is only an estimate, simply decrease η . A smaller η still yields strong disentanglement (no wrap), at the cost of reduced dynamic range, which is usually acceptable.

Intuitively, η directly scales the maximum realized rotation. When in doubt, pick a smaller η ; this preserves the equality (no modulo) in equation 20.

Having established strong disentanglement, we can now read out the latent coordinate on each style plane by a simple geometric post-processing of the feature angles.

On a non-cyclic factor, the feature angle on the circle is already an affine function of the latent, but only defined modulo 2π . For each style subspace, the data only occupies one continuous arc on the circle and there is a single "big gap" where the angles wrap around from the end back to the beginning. The construction in Corollary 2 simply detects this largest gap, takes the midpoint of the covered arc as a new reference direction, and then re-centers all angles around it. With this re-centering, the wrap disappears and the principal angle becomes a true linear function of the latent, up to a shift by the mid-point.

For cyclic factors, one full latent period must correspond to exactly one full 2π turn on the circle. Strong disentanglement forces the rotation angle to grow linearly with the latent, so "one period \Rightarrow one full turn" immediately pin down the slope. The remaining freedom is just a global in-plane rotation that is fixed per cluster, so within each cluster the feature on that style plane is exactly a phase-shifted sine-cosine trace of the latent variable which preserve the cyclic structure, as shown in Corollary 3.

Corollary 2 (Principal-angle post-processing when the largest gap is the endpoint wrap gap). Fix the i -th style plane and an \mathcal{A} -connected cluster c_k with non-cyclic latent support $z^{(i)} \in [a_k, b_k]$ (the range may depend on the cluster c_k). For each $x = \varphi(z)$ let $\tilde{f}_{\text{sty}}^{(i)}(x) = (u_i(x), v_i(x)) \in \mathbb{S}^1$ and

$$\theta(x) := \text{atan2}(v_i(x), u_i(x)) \in (-\pi, \pi].$$

Assume the plane-wise strong disentanglement condition (or, equivalently, weak edge-wise equivariance + no-wrap span) so that there exist $c_0^{(i)} \neq 0$ and $\beta^{(i)} \in \mathbb{R}/2\pi\mathbb{Z}$ with

$$\theta(x) \equiv c_0^{(i)} z^{(i)} + \beta^{(i)} \pmod{2\pi}. \quad (26)$$

Let the principal angles be sorted $-\pi < \theta_1 < \dots < \theta_n \leq \pi$, and define gaps $g_k := \theta_{k+1} - \theta_k$ for $k = 1, \dots, n-1$ and $g_n := \theta_1 + 2\pi - \theta_n$. Let $k^* = \arg \max_k g_k$ and assume the largest gap g_{k^*} is the wrap-around gap between the endpoint latents a_k and b_k (endpoint wrap-gap).

Center of the covered arc (no explicit wrapping). Let

$$\begin{aligned} \theta_L &:= \theta_{k^*}, & \theta_R &:= \theta_{k^*+1} \text{ (with } \theta_{n+1} := \theta_1 + 2\pi), \\ e_L &= (\cos \theta_L, \sin \theta_L), & e_R &= (\cos \theta_R, \sin \theta_R), & m &:= e_L + e_R. \end{aligned}$$

Set

$$\alpha := \begin{cases} \operatorname{atan2}(m_y, m_x), & \text{if } g_{k^*} > \pi \text{ (} L := 2\pi - g_{k^*} < \pi), \\ \operatorname{atan2}(-m_y, -m_x), & \text{if } g_{k^*} \leq \pi \text{ (} L := 2\pi - g_{k^*} \geq \pi), \end{cases} \quad (27)$$

where L is the length of the covered data arc (the complement of the largest gap). Finally, define the 1-D feature

$$\widehat{f}^{(i)}(x) := \operatorname{atan2}(\sin(\theta(x) - \alpha), \cos(\theta(x) - \alpha)) \in (-\pi, \pi]. \quad (28)$$

Conclusion (mid-latent linearity). For all $x = \varphi(z)$,

$$\widehat{f}^{(i)}(x) = c_0^{(i)} \left(z^{(i)} - \frac{a_k + b_k}{2} \right),$$

i.e., an exact, continuous affine readout centered at the mid-latent $(a_k + b_k)/2$; if $c_0^{(i)} > 0$ it is strictly increasing, and if $c_0^{(i)} < 0$ it is strictly decreasing.

Proof. For any $\varphi(z_0) = x_0 \sim_{\mathcal{A}} x = \varphi(z)$. Since unit vectors $a = (\cos \alpha, \sin \alpha)$ and $b = (\cos \beta, \sin \beta)$, $\operatorname{atan2}(a \times b, a \cdot b) \equiv \beta - \alpha \pmod{2\pi}$, applying this to $a = \tilde{f}_{\text{sty}}^{(i)}(x_0)$, $b = \tilde{f}_{\text{sty}}^{(i)}(x)$ gives

$$\theta(x) - \theta(x_0) \equiv \theta_i^{\text{pr}}(\tilde{f}(x); \tilde{f}(x_0)) \pmod{2\pi}.$$

Since $\theta_i^{\text{dir}} \equiv \theta_i^{\text{pr}} \pmod{2\pi}$ by definition of θ_i^{dir} , strong disentanglement yields for every connected pair:

$$\theta(x) - \theta(x_0) \equiv c_0^{(i)} (z^{(i)} - z_0^{(i)}) \pmod{2\pi}.$$

Set $\beta^{(i)} := \theta(x_0) - c_0^{(i)} z_0^{(i)} \in \mathbb{R}/2\pi\mathbb{Z}$ to obtain equation 26.

(1) *Largest gap \Rightarrow covered arc and its endpoints.* By construction the n gaps $\{g_k\}$ partition the circle and sum to 2π . Removing the *largest* gap g_{k^*} leaves the *shortest* circular arc that covers all data; its endpoints are precisely the two unit directions with principal angles θ_R and θ_L adjacent across g_{k^*} . By the endpoint wrap-gap assumption, these endpoints are the samples at the latent endpoints $z^{(i)} = a$ and $z^{(i)} = b$, so the covered-arc length is $L = 2\pi - g_{k^*} \in (0, 2\pi)$.

(2) *Center via the two endpoint vectors.* For two unit vectors at angles ϕ_1, ϕ_2 , the vector sum has direction

$$\arg(e^{i\phi_1} + e^{i\phi_2}) = \frac{\phi_1 + \phi_2}{2}$$

and points to the midpoint of the *shorter* of the two arcs between ϕ_1 and ϕ_2 . In our notation, $m = e_L + e_R$ points to the midpoint of the shorter arc between θ_L and θ_R . If $g_{k^*} > \pi$ then the gap is longer than the data arc ($L < \pi$), so the shorter arc is the data arc and $\arg(m)$ is already the covered-arc midpoint. If $g_{k^*} \leq \pi$ then the gap is the shorter arc (and the data arc has length $L \geq \pi$), so the covered-arc midpoint is the *antipode* of $\arg(m)$; this is achieved by $\arg(-m)$. This is exactly the definition equation 27 of α , which lies in $(-\pi, \pi]$ by the properties of $\operatorname{atan2}$.

(3) *Signed offsets about α are single-branch differences.* All principal angles $\theta(x)$ lie in the covered arc of length $L < 2\pi$ centered at α , hence for every sample $|\theta(x) - \alpha| \leq L/2 < \pi$. Therefore, the signed principal difference computed by equation 28 equals the ordinary difference:

$$\widehat{f}^{(i)}(x) = \theta(x) - \alpha \quad \forall x.$$

(4) *Mid-latent identity for α and exact linearity.* By the modular–affine law equation 26, there exist integers m_a, m_b such that the *unwrapped* endpoint angles along the covered arc are

$$\phi(a_k) = c_0^{(i)} a_k + \beta^{(i)} + 2\pi m_a, \quad \phi(b_k) = c_0^{(i)} b_k + \beta^{(i)} + 2\pi m_b,$$

with $\phi(b) - \phi(a) = L \in (0, 2\pi)$. The covered-arc midpoint is then

$$\alpha \equiv \frac{\phi(a_k) + \phi(b_k)}{2} \equiv c_0^{(i)} \frac{a_k + b_k}{2} + \beta^{(i)} \pmod{2\pi}.$$

Because every $\theta(x)$ lies on the single branch centered at α (Step 3), the congruence equation 26 lifts to the equality $\theta(x) = c_0^{(i)} z^{(i)} + \beta^{(i)}$ on this branch. Subtracting $\alpha \equiv c_0^{(i)} \frac{a_k + b_k}{2} + \beta^{(i)}$ yields

$$\tilde{f}^{(i)}(x) = \theta(x) - \alpha = c_0^{(i)} \left(z^{(i)} - \frac{a_k + b_k}{2} \right).$$

Continuity follows since no sample sits on the cut, and Step 3 used signed principal differences. \square

Corollary 3 (Single-period cyclic factor determines the slope (clusterwise phase)). *Fix a style plane j and an \mathcal{A} -connected cluster c_k whose latent support on this plane is intrinsically cyclic with period $p > 0$. Suppose there exists an interval $I = [a, a + p)$ (a single latent period) across which $z^{(j)}$ varies on the data support of c_k . Assume the plane-wise strong disentanglement. Then:*

$$|c_0^{(j)}| = \frac{2\pi}{p}, \quad \text{and} \quad \exists A_{j,k} \in SO(2) \text{ (independent of } z) \text{ s.t. for all } z^{(j)} \in I,$$

$$\tilde{f}_{\text{sty}}^{(j)}(\varphi(z)) = A_{j,k} \begin{pmatrix} \cos\left(\frac{2\pi}{p} z^{(j)}\right) \\ \sin\left(\frac{2\pi}{p} z^{(j)}\right) \end{pmatrix}.$$

In particular, within each cluster c_k the representation on plane j is exactly periodic with period p and linear in $z^{(j)}$ modulo a clusterwise fixed in-plane rotation $A_{j,k}$.

Proof. Fix a cluster k , pick any $z_0 \sim_{\mathcal{A}} z$ and by the perfect equivariance condition,

$$\tilde{f}_{\text{sty}}^{(j)}(\varphi(z)) = R(c_0^{(j)}(z^{(j)} - z_0^{(j)})) \tilde{f}_{\text{sty}}^{(j)}(\varphi(z_0)), \quad \forall z \in \mathcal{Z}_{\text{train}}. \quad (29)$$

As $z^{(j)}$ moves from a to $a + p$, the feature’s rotation angle changes by $c_0^{(j)}((a + p) - a) = c_0^{(j)}p$. Strong disentanglement makes $\tilde{f}_{\text{sty}}^{(j)}(\varphi(z))$ traverse the circle exactly once, forcing $|c_0^{(j)}p| = 2\pi$, hence $|c_0^{(j)}| = 2\pi/p$.

Let $\phi_0 \in \mathbb{R}$ be the phase of $\tilde{f}_{\text{sty}}^{(j)}(\varphi(z_0))$, i.e. $\tilde{f}_{\text{sty}}^{(j)}(\varphi(z_0)) = (\cos \phi_0, \sin \phi_0)$. Plug $c_0^{(j)}$ in equation 29:

$$\tilde{f}_{\text{sty}}^{(j)}(\varphi(z)) = R\left(\frac{2\pi}{p}(z^{(j)} - z_0^{(j)})\right) R(\phi_0) e_1 = R\left(\frac{2\pi}{p} z^{(j)} + \underbrace{\phi_0 - \frac{2\pi}{p} z_0}_{=: \theta_j}\right) e_1,$$

where we used $R(\alpha)R(\beta) = R(\alpha + \beta)$ and commutativity of planar rotations.

Claim. θ_j is independent of the choice of z_0 .

Pick another reference $z_1 \sim_{\mathcal{A}} z$. Using the equivariance condition, we have

$$\tilde{f}_{\text{sty}}^{(j)}(\varphi(z_1)) = R\left(\frac{2\pi}{p}(z_1^{(j)} - z_0^{(j)})\right) \tilde{f}_{\text{sty}}^{(j)}(\varphi(z_0))$$

Writing $\tilde{f}_{\text{sty}}^{(j)}(\varphi(z_1)) = (\cos \phi_k, \sin \phi_k)$, this implies $\phi_1 \equiv \phi_0 + \frac{2\pi}{p}(z_1^{(j)} - z_0^{(j)}) \pmod{2\pi}$. Therefore

$$\phi_1 - \frac{2\pi}{p} z_1^{(j)} \equiv \left(\phi_0 + \frac{2\pi}{p}(z_1^{(j)} - z_0^{(j)})\right) - \frac{2\pi}{p} z_1^{(j)} = \phi_0 - \frac{2\pi}{p} z_0^{(j)} \pmod{2\pi},$$

and since $R(\cdot)$ is 2π -periodic, $R(\phi_1 - \frac{2\pi}{p} z_1) = R(\phi_0 - \frac{2\pi}{p} z_0)$. Hence, prove the claim. Define the fixed rotation

$$A_{j,k} := R(\theta_j) \in SO(2).$$

1674 Then

$$1675 \tilde{f}_{\text{sty}}^{(j)}(\varphi(z)) = A_{j,k} R\left(\frac{2\pi}{p} z^{(j)}\right) e_1 = A_{j,k} \begin{pmatrix} \cos\left(\frac{2\pi}{p} z^{(j)}\right) \\ \sin\left(\frac{2\pi}{p} z^{(j)}\right) \end{pmatrix}, \quad \forall z \in \mathcal{Z}_{\text{train}}, \quad (30)$$

1676 which is the claimed representation with plane rotation $A_{j,k}$ independent of reference z_0 . \square

1680 **Remark 4** (Global alignment after post-processing). *Non-cyclic case.* If for plane i all \mathcal{A} -connected clusters share the same midpoint $(a_i + b_i)/2$, then the principal-angle post-processing of Cor. 2 yields

$$1681 \hat{f}^{(i)}(x) = c_0^{(i)} \left(z^{(i)} - \frac{a_i + b_i}{2} \right) \quad \text{for every cluster,}$$

1682 so the 1-D readout is already aligned globally across clusters.

1686 *Cyclic case.* For intrinsically cyclic planes, there is no lossless 1-D readout; features remain 2-D on S^1 and differ across clusters by a constant phase $A_{j,k} \in SO(2)$. Consequently, global comparison requires a clusterwise phase alignment (e.g., rotate each cluster by $A_{j,k}^{-1}$ to a chosen reference gauge); after this rotation, all clusters live in the same 2-D gauge up to one shared global in-plane rotation.

1689 In the last two propositions, we only prove the continuous case, which can be easily extended to the discrete case.

1694 C.5 PROOF OF PROPOSITION 4

1696 **Proposition 8** (Axis-separable linearity extension). *Under the assumption of Corollary 1. Moreover, except single-style augmentations $\mathcal{A}_{\text{single}} = \{a_i(t_i)\}_{i=1}^{d_{\text{sty}}}$, there are $a_{\text{as}}(t) \in \mathcal{A} \setminus \mathcal{A}_{\text{single}}$ such that:*

$$1699 a_{\text{as}}(t_{\text{as}}, \cdot) = a_{i_m}(t_m) \circ \cdots \circ a_{i_1}(t_1, \cdot),$$

1700 (Here “ \circ ” means usual function composition: $(a \circ b)(x) = a(b(x))$.)

1702 Suppose for all $t_{\text{as}} \in \mathcal{T}_{\text{as}}$ there exists z_{as} and a finite chain

$$1703 z_0 := z_{\text{as}}, z_1, \dots, z_{n-1} \in \mathcal{Z}_{\text{train}}, \quad z_n := a_{\text{as}}(t_{\text{as}}, z_{\text{as}}) \quad (z_n \text{ not required in } \mathcal{Z}_{\text{train}}),$$

1705 such that for each edge (z_j, z_{j+1}) with $j = 0, \dots, n-2$, there exist augmentations $u_j^+(t_j^+), u_j^-(t_j^-)$ such that

$$1706 \varphi^{-1}(u_j^+(t_j^+, \varphi(z_j))) = \varphi^{-1}(u_j^-(t_j^-, \varphi(z_{j+1}))) =: z_j^* \in \mathcal{Z},$$

1708 realized displacements $m_j^+ := u_j^+(t_j^+)$, $m_j^- := u_j^-(t_j^-) \in \mathbb{R}^{d_{\text{sty}}}$ satisfying $m_j^{+,(i)}, m_j^{-,(i)} \in [-M_i, M_i]$ for all i . For the last edge (z_{n-1}, z_n) , there exists a one-sided augmentation $u_{n-1}^+(t_{n-1}^+)$ with

$$1710 \varphi(z_n) = u_{n-1}^+(t_{n-1}^+, \varphi(z_{n-1})), \quad m_{n-1}^+ := u_{n-1}^+(t_{n-1}^+), \quad m_{n-1}^{+,(i)} \in [-M_i, M_i].$$

1714 Then $a_{\text{as}}(t_{\text{as}})$ is axis-separable with net latent displacement

$$1715 v := \sum_{j=0}^{n-2} (m_j^- - m_j^+) + m_{n-1}^+$$

1716 and $g_{\text{as}} : \mathcal{T}_{\text{as}} \rightarrow [-\pi, \pi]^{d_{\text{sty}}}$ is linear in the per-axis displacement:

$$1717 g_{\text{as}}(t_{\text{as}}) \equiv c_0 \odot v \quad (\text{mod } 2\pi \text{ per axis})$$

1722 If in addition, $v^{(i)} \in [-M_i, M_i]$, then

$$1723 g_{\text{as}}(t_{\text{as}}) = c_0 \odot v$$

1724 *Proof.* Write $x_j := \varphi(z_j)$. For each interior edge, perfect equivariance on the two routes to z_j^* gives $\mathcal{R}(c_0 \odot m_j^+) \tilde{f}(x_j) = \tilde{f}(x_j^*) = \mathcal{R}(c_0 \odot m_j^-) \tilde{f}(x_{j+1})$, hence $\tilde{f}(x_{j+1}) = \mathcal{R}_{c_0 \odot (m_j^- - m_j^+)} \tilde{f}(x_j)$.

For the terminal edge, one-sided equivariance yields $\tilde{f}(x_n) = \mathcal{R}_{c_0 \odot m_{n-1}^+} \tilde{f}(x_{n-1})$. Multiplying all steps and using those rotations on distinct planes commute and add in-plane,

$$\tilde{f}(x_n) = \mathcal{R}(c_0 \odot (\sum_{j=0}^{n-2} (m_j^- - m_j^+) + m_{n-1}^+)) \tilde{f}(x_0) = \mathcal{R}_{c_0 \odot v} \tilde{f}(x_0),$$

Since the output of g_{as} is bounded in $[-\pi, \pi]$, we can only have $g_{\text{as}}(t_{\text{as}}) \equiv c_0 \odot v \pmod{2\pi}$ per axis).

If in addition, we have $v^{(i)} \in [-M_i, M_i]$, since each g_i output with bound $[-\pi, \pi]$, we can have $g_{\text{as}}(t_{\text{as}}) = c_0 \odot v$.

□

C.6 PROOF OF PROPOSITION 5

Proposition 9 (Few-latent anchor, gated composite head, and linearity). *Let $\mathcal{A}_{\text{single}} = \{a_i(t_i)\}_{i=1}^{d_{\text{sty}}-1}$ be single-style augmentations with z -independent per-axis displacements $\Delta_{a_i}(t_i, z) = q_{a_i}(t_i) e_i$ and $|q_{a_i}(t_i)| \leq M_i$. Let $a_{\text{cmp}}(t) \in \mathcal{A} \setminus \mathcal{A}_{\text{single}}$ be a (possibly) multi-style augmentation.*

Assume:

(A1) **Pure at an anchor fiber.** *There is $J \subseteq \{1, \dots, d_{\text{sty}} - 1\}$ and \bar{z}_J such that on the fiber $F_J(\bar{z}_J) := \{z : \pi_J(z) = \bar{z}_J\}$,*

$$\Delta_{a_{\text{cmp}}}(t, z) = q_{a_{\text{cmp}}}(t) e_{d_{\text{sty}}}, \quad |q_{a_{\text{cmp}}}(t)| \leq M_{d_{\text{sty}}}.$$

(A2) **Axis-definable off the anchor.** *For any (t, z) there exists an \mathcal{A} -chain (as in Lemma 8) from z to $a_{\text{cmp}}(t, z)$ comprised of realized steps m_j^\pm (each coordinatewise bounded by M_i) with net*

$$v(t, z) := \sum_{j=0}^{n-2} (m_j^- - m_j^+) + m_{n-1}^+ \in \mathbb{R}^{d_{\text{sty}}}.$$

Gated composite head. *Define the selector $S : \mathbb{R}^{|J|} \rightarrow \{0, 1\}^{d_{\text{sty}}}$ by*

$$S(\pi_J(z)) = \begin{cases} e_{d_{\text{sty}}}, & \pi_J(z) = \bar{z}_J \quad (\text{on the anchor fiber}), \\ \mathbf{1}_{d_{\text{sty}}}, & \pi_J(z) \neq \bar{z}_J \quad (\text{off fiber}), \end{cases}$$

and take

$$g_{\text{cmp}} : \mathbb{R}^{|J|} \times \mathcal{T}_{a_{\text{cmp}}} \rightarrow [-\pi, \pi]^{d_{\text{sty}}}, \quad g_{\text{cmp}}(\pi_J(z), t) \in S(\pi_J(z)) \odot [-\pi, \pi]^{d_{\text{sty}}}.$$

Conclusions. *Assume perfect equivariance (Lemma 2), nondegeneracy $g_{\text{cmp}} \not\equiv 0$ on non-identity displacement, and Cor. 1 for single-axis and pure-located cases. Then:*

(C1) **On the anchor (one channel open).** *For any $z \in F_J(\bar{z}_J)$ there exists $c_0^{(d_{\text{sty}})} \neq 0$ s.t.*

$$g_{\text{cmp}}(\bar{z}_J, t) = c_0^{(d_{\text{sty}})} q_{a_{\text{cmp}}}(t).$$

(C2) **Off the anchor (all channels open).** *For any $z \notin F_J(\bar{z}_J)$,*

$$g_{\text{cmp}}(\pi_J(z), t) \equiv c_0 \odot v(t, z)$$

for some constant slope vector $c_0 \in \mathbb{R}^{d_{\text{sty}}}$ determined by the single-axis calibrations.

Proof. Step 1: Equivariance to in-plane rotations. By perfect equivariance, for any augmentation that effects a net latent displacement $w \in \mathbb{R}^{d_{\text{sty}}}$, the aligned style features transform as

$$\tilde{f}(a(\cdot, \varphi(z))) = \mathcal{R}(\theta(w)) \tilde{f}(\varphi(z)), \quad \text{with } \theta(w) = c_0 \odot w,$$

where c_0 is the (cluster/plane-wise) slope vector supplied by Cor. 1 from the single-axis identifications.

Step 2: Anchor purity \Rightarrow scalar linearity. On $F_J(\bar{z}_J)$, assumption (A1) gives $\Delta_{a_{\text{cmp}}}(t, z) = q_{a_{\text{cmp}}}(t)e_{d_{\text{sty}}}$. Since the selector $S(\bar{z}_J) = e_{d_{\text{sty}}}$ opens only channel d_{sty} , g_{cmp} depends only on t but not on other coordinates. By the single-axis linearity (Cor. 1) applied on the d_{sty} -plane,

$$g_{\text{cmp}}(\bar{z}_J, t) = c_0^{(d_{\text{sty}})} q_{a_{\text{cmp}}}(t),$$

with $c_0^{(d_{\text{sty}})} \neq 0$ because g_{cmp} is nontrivial on non-identity displacements. This proves (C1).

Step 3: Axis-definability off fiber \Rightarrow vector linearity. Off the anchor, (A2) provides an \mathcal{A} -chain whose realized net is $v(t, z)$. Equivariance composes along the chain (Lemma 8): each single-axis step contributes additively to the angle in its plane, so the total in-plane rotation equals $\theta(v(t, z)) = c_0 \odot v(t, z)$. Because $S(\pi_J(z)) = \mathbf{1}_{d_{\text{sty}}}$, the head exposes all coordinates and thus reads out

$$g_{\text{cmp}}(\pi_J(z), t) \equiv c_0 \odot v(t, z).$$

This yields (C2).

Step 4: No-wrap regime under calibration. If each coordinate is calibrated so that no coordinate wraps and the equality holds in $\mathbb{R}^{d_{\text{sty}}}$, $g_{\text{cmp}} = c_0 \odot v(t, z)$. \square

Remark 5 (Training with “pure indicators”, fixed gating code). *We do not observe latent z . We mark a subset $\mathcal{X}_{\text{pure}} \subset \mathcal{X}_{\text{train}}$ on which a_{cmp} is empirically pure on the d_{sty} -plane and assign a fixed code to all such points:*

$$h(x) \equiv \mathbf{0} \in \mathbb{R}^{|J|} \quad \text{for all } x \in \mathcal{X}_{\text{pure}} \quad (\text{no learning of } h).$$

The gate uses $h(x)$ in place of $\pi_J(z)$: on $\mathcal{X}_{\text{pure}}$, set $S(h(x)) = e_{d_{\text{sty}}}$ (only channel d_{sty} open); off that set, open all channels, e.g. $S(h(x)) = \mathbf{1}_{d_{\text{sty}}}$.

A practical loss is

$$\mathcal{L}_{\text{gate}} := \mathbb{E}_{x \in \mathcal{X}_{\text{pure}}, t} \sum_{i \neq d_{\text{sty}}} (g_{\text{cmp}}^{(i)}(h(x), t))^2,$$

which closes all non- d_{sty} channels on pure points; off-fiber behavior is learned via the \mathcal{A} -chain supervision.

Why fixing h on pure points is necessary. *If $h(x)$ were allowed to vary across $x \in \mathcal{X}_{\text{pure}}$, then equivariance alone would not guarantee a single angle $\theta^{(d_{\text{sty}})}$ for the same augmentation t across pure samples: the head could implement a sample-dependent reparametrization $\theta^{(d_{\text{sty}})}(t, h(x))$, still satisfying equivariance but breaking identifiability of the common slope on the d_{sty} -plane. Fixing $h(x) \equiv \mathbf{0}$ eliminates this degree of freedom and enforces a unique readout for the same augmentation.*

1836 **Algorithm 1** Cluster-Dependent Rotational Equivariance (with post-processing)

1837 1: **Input** Dataset $\{x\}$, augmentations $\mathcal{A} = \{a_i\}_{i=1}^m$, number of clusters k , transform nets $G =$
1838 $\{g_i\}_{i=1}^m$, feature extractor f , hyperparams $\{\lambda_{\text{theta}}, \lambda_{\text{radius}}\}$, subspace radius ω , epochs T_1, T_2

1839 **Stage 1 (optional contrastive pretrain; only if $k > 1$)**

1840 2: **for** $t = 1$ to T_1 **do**

1841 3: Sample $\{(x_j, x_j^\dagger = a(x_j))\}$ with $a \in \mathcal{A}$

1842 4: Compute $\mathcal{L}_{\text{InfoNCE}}$

1843 5: Update f by $\nabla \mathcal{L}_{\text{InfoNCE}}$

1844 6: **end for**

1845 **Stage 2 (Cluster-Dependent Rotational Equivariance)**

1846 7: Initialize centers $\{c_\ell\}$ by spherical k -means on $u(x) := f(x)/\|f(x)\|_2$;

1847 8: Set $r_\ell \leftarrow \text{normalize}(c_\ell)$; Constrain range of g_i by Eq. (13)

1848 9: **for** $t = 1$ to T_2 **do**

1849 10: Sample $\{(x_j, x_j^\dagger = a(t, x_j))\}$ with $a(t) \in \mathcal{A}$

1850 11: $u_j \leftarrow f(x_j)/\|f(x_j)\|_2$, $u_j^\dagger \leftarrow f(x_j^\dagger)/\|f(x_j^\dagger)\|_2$

1851 12: Assign $c(x_j) = \arg \max_i u_j^\top r_i$ (same for x_j^\dagger)

1852 13: Build Householder $H_{c(x)}$ that maps the rotation center $r_{c(x)}$ pole to the north pole

1853 14: $\theta \leftarrow G(t)$

1854 15: Build block-diag matrix $R(\theta)$

1855 16: Compute $\mathcal{L}_{\text{equiv}}, \mathcal{L}_{\text{radius}}, \mathcal{L}_{\text{theta}}$

1856 17: Update $f, \{g_i\}$ by $\nabla(\mathcal{L}_{\text{equiv}} + \lambda_{\text{radius}}\mathcal{L}_{\text{radius}} + \lambda_{\text{theta}}\mathcal{L}_{\text{theta}})$

1857 18: **end for**

1858 **Optional: post-processing and final representation**

1859 19: **for** each cluster i and style plane j **do**

1860 20: Collect $\theta(x) \leftarrow \text{atan2}([H_i f(x)]_{j,y}, [H_i f(x)]_{j,x})$ for $x \in S_i$

1861 21: Sort $\{\theta(x)\}$ to $\theta_1 < \dots < \theta_n$; $g_k \leftarrow \theta_{k+1} - \theta_k$ ($k = 1:n-1$), $g_n \leftarrow \theta_1 + 2\pi - \theta_n$

1862 22: $k^* \leftarrow \arg \max_k g_k$; $\theta_L \leftarrow \theta_{k^*}, \theta_R \leftarrow \theta_{k^*+1}$ (wrap)

1863 23: $m \leftarrow (\cos \theta_L, \sin \theta_L) + (\cos \theta_R, \sin \theta_R)$; $L \leftarrow 2\pi - g_{k^*}$

1864 24: $\alpha_{i,j} \leftarrow \begin{cases} \text{atan2}(m_y, m_x), & L < \pi \\ \text{atan2}(-m_y, -m_x), & L \geq \pi \end{cases}$

1865 25: **end for**

1866 26: **for** each x with cluster $i = c(x)$ **do**

1867 27: **for** each plane j **do**

1868 28: **if** plane j is non-cyclic **then**

1869 29: $\tilde{f}^{(j)}(x) \leftarrow \text{atan2}(\sin(\theta^{(j)}(x) - \alpha_{i,j}), \cos(\theta^{(j)}(x) - \alpha_{i,j}))$ (scalar)

1870 30: **else**

1871 31: $\tilde{f}^{(j)}(x) \leftarrow A_{i,j} [H_i f(x)]_j \in \mathbb{R}^2$ (with $A_{i,j}$ fixed per cluster/plane)

1872 32: **end if**

1873 33: **end for**

1874 34: $\hat{f}_{\text{style}}(x) \leftarrow \text{concat}(\{\tilde{f}^{(j)}(x)\}_{j \in \text{non-cyc}}, \{\tilde{f}^{(j)}(x)\}_{j \in \text{cyc}})$

1875 35: **Final rep:** $\hat{r}(x) \leftarrow \text{concat}(\hat{f}_{\text{style}}(x), e_i)$ (append one-hot cluster code $e_i \in \{0, 1\}^k$)

1876 36: **end for**

Table 2: Parameter ranges used in our setup.

Parameter	Minimum value	Maximum value
Object position X	-0.3	0.3
Object position Y	-0.3	0.3
Object position Z	-0.3	0.3
Object rotation ψ	-0.5	0.5
Floor hue	0.3	0.7

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

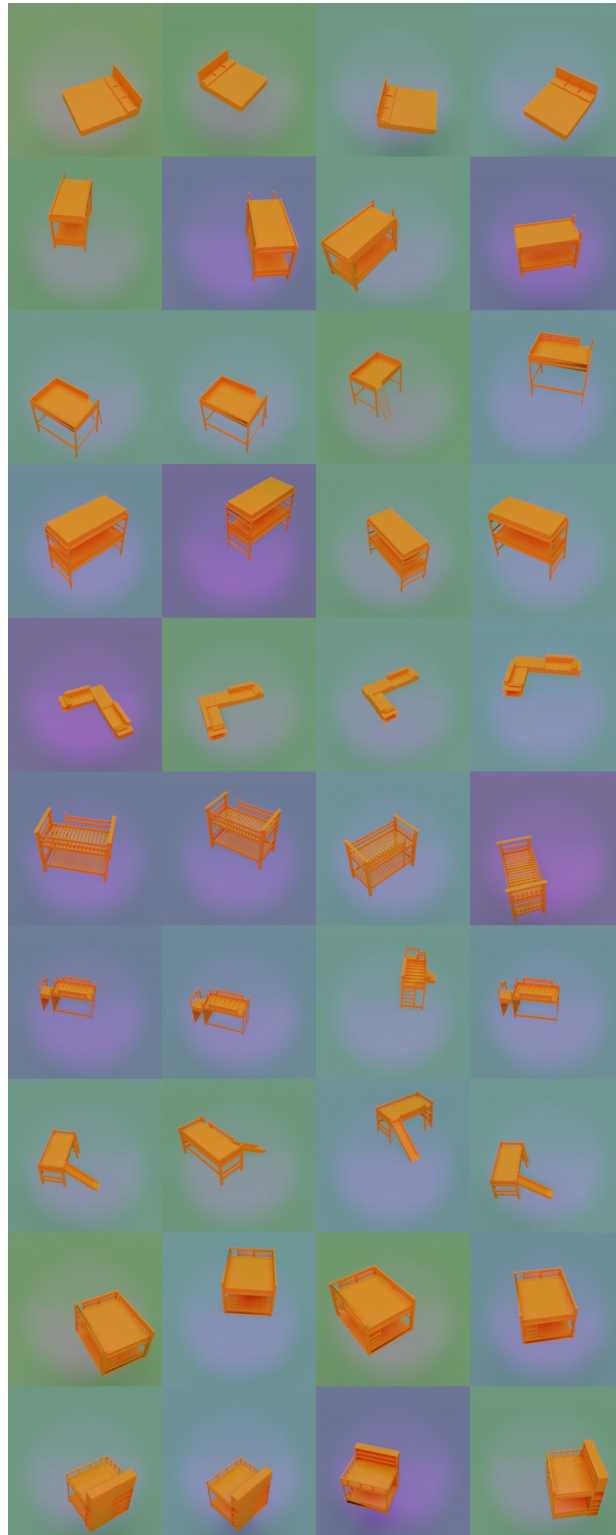


Figure 5: Samples from created 3DIEBench.

D EXPERIMENT DETAILS

D.1 DATASET

Shape3D (Kim & Mnih, 2018) This dataset provides 5 style latents (floor hue, wall hue, object hue, object scale, azimuth) and 4 object shapes as the class latent. For each style latent, we keep the middle half of its range and include all 4 classes, yielding a total of 16,000 images.

MPI3D (Gondal et al., 2019) This dataset contains 6 style latents (object color, object size, camera height, background color, horizontal rotation, vertical rotation) and 6 object shapes as the class latent. We keep the first four style latents in full. For the two rotations, we subsample the angles to $\{0, 4, 8, 12\}$ (uniformly spaced indices), resulting in a total of 10,384 images.

3DIdent (Zimmermann et al., 2021) This dataset has a single class and 10 continuous style factors: position (x, y, z) ; rotation (ϕ, θ, ψ) ; spotlight position; object hue; ground hue; and spotlight hue. We randomly sample 50,000 images for training.

3DIEBench (Garrido et al., 2023) 3DIEBench is derived from a ShapeNetCore subset (sourced from 3D Warehouse). To generate the non-symmetric images, we use 10 chair models as classes and vary 5 parameters using Blender to render images uniformly over their ranges (times π in the generating object rotation ψ). We generate 20,000 images at a resolution of $224 * 224$. Parameter ranges are listed in Table 2, and example renderings are shown in Figure 5.

Synthetic dataset The synthetic dataset used in Section 4.3 is designed so that we can precisely control the underlying factors of variation and directly verify our theoretical results. It consists of 96×96 grayscale images of three distinct uppercase letters, “R,” “E,” and “D.” For each letter, we generate samples by applying three controlled augmentations: horizontal translation (x-shift), vertical translation (y-shift), and in-plane rotation about the image center.

D.2 HARDWARE

All experiments were performed on 2 NVIDIA Tesla V100 GPUs with 32GB accelerator RAM for a single training run. All experiments use the PyTorch deep learning framework (Paszke et al., 2019).

D.3 TRAINING

Unless otherwise stated, all methods are trained for 500 epochs with a ResNet-18 backbone and a three-layer MLP projection head (hidden sizes 2048-2048-2048, PReLU activations). We use a batch size of 500 and Adam ($\text{lr} = 10^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$). Feature dimension is set to be $2 * \# \text{ latents} + 1$, where 1 is the extra dimension as stated in main content. Method-specific hyperparameters are detailed below.

SimCLR (Chen et al., 2020a). We set the InfoNCE temperature to $\tau = 0.5$ for all experiments.

CARE (Gupta et al., 2023). We use $\tau = 0.5$ in the InfoNCE loss and $\lambda_{\text{equi}} = 10$. The number of equivariance chunks is set to half the batch size ($B/2$).

EquiMod (Devillers & Lefort, 2022). We use $\tau = 0.5$ for InfoNCE loss and $\lambda_{\text{equimod}} = 1$. The predictor head mirrors the projection head architecture.

IP-IRM (Wang et al., 2021). We use $\tau = 0.5$ for InfoNCE loss and $\lambda_1 = \lambda_2 = 0.2$ and change each step per 30 epochs. For the updating of partitions, the training epochs are set to be 50.

Eastwood et al. (2023). We use $\tau = 0.5$ for all the losses and for each subspaces we set the dimension = 2.

Ours. We adopt a two-stage schedule. Stage 1: SimCLR pretraining for 100 epochs with $\tau = 0.5$ (skipped on **3DIdent**). Stage 2: we train the transformation networks $G = \{g_i\}_i$ as three-layer MLPs (hidden sizes 2048–2048–2048, PReLU + LayerNorm) for the remaining 400 epochs (or 500 epochs on **3DIdent** when Stage 1 is skipped). Subspace radius ω is set to be 0.05 and extra dimension is set to be 1 for all the datasets. We simply set loss weights $\lambda_{\text{radius}} = \lambda_{\text{theta}} = 1$ for all four datasets.

D.4 EVALUATION

We evaluate with **DCI** (Eastwood & Williams, 2018), which summarizes three properties of a representation: *disentanglement*, *completeness*, and *informativeness*. Following Eastwood & Williams (2018), we train a supervised predictor from the learned features to the ground-truth factors and extract an *importance matrix* $R \in \mathbb{R}_{\geq 0}^{d \times k}$, where R_{ij} measures the contribution of feature dimension i to predicting factor j . In our implementation, we use either a Linear regressor or a Random Forest regressor: for Linear, we take $R_{ij} = |w_{ij}|$ (absolute regression weights), and for Random Forest, we use the mean decrease in impurity (Gini importance) (Breiman, 2001).

Let \tilde{R} be R normalized to probabilities: row-wise for disentanglement and column-wise for completeness. Disentanglement and completeness are then computed as one minus the normalized entropy of these distributions and aggregated across dimensions/factors with standard DCI weighting; informativeness is the prediction performance of the supervised model (we report R^2 / error as appropriate), exactly as in Eastwood & Williams (2018).

Because our post-processed features include a one-hot class indicator, we treat that block as a *single* semantic factor. Before normalization, we collapse the corresponding columns by summing their importances, so that the class indicator contributes as one factor in R .

E ABLATION

We ablate both the training protocol and the loss terms. (1) Without the SimCLR pretraining stage—i.e., if we initialize rotation centers by running k -means on randomly initialized features—the learned features degrade markedly. (2) Even with proper SimCLR pretraining, removing either the radius penalty $\mathcal{L}_{\text{radius}}$ or the angle penalty $\mathcal{L}_{\text{theta}}$ leads to feature collapse (features concentrate to a single point on the sphere). Quantitative and qualitative evidence is provided in Table 3 and Figure 6.

Table 3: Per-dimension Informativeness using linear regression on **Shapes3D**. Columns list ground-truth factors

Method	Floor Hue	Wall Hue	Object Hue	Scale	Orientation	Class
w/o pretraining	0.0648	0.2712	0.0177	0.0172	0.0010	0.0034
w/o $\mathcal{L}_{\text{radius}}$	0.0000	0.0000	0.0000	0.0000	0.0000	1.0000
w/o $\mathcal{L}_{\text{theta}}$	0.0000	0.0000	0.0000	0.0000	0.0000	1.0000
CD-RED (after post-proc.)	0.9995	0.9995	0.9994	0.9984	0.9994	1.0000

Uniform rotation centers (optional). To ensure maximally separated clusters, we may further add a repulsion term over rotation centers

$$\mathcal{L}_{\text{uniform}} = -\mathbb{E}_{i \neq j} \|r_i - r_j\|_2^2, \quad (31)$$

which, when minimized jointly with our main losses, promotes a uniform configuration on the sphere. So each rotation center is initialized as the cluster centers right after the clustering. This stabilizes updates while allowing $\mathcal{L}_{\text{uniform}}$ to shape inter-center geometry.

Under the same settings as Figure 4, we report the minimum pairwise center distance $d_{\text{min}} = \min_{i \neq j} \|r_i - r_j\|_2$. With learnable rotation centers (i.e., $\mathcal{L}_{\text{uniform}}$ on), the configuration approaches the Tammes-optimal spacing: for $k = 3$,

$$d_{\text{min}} = \sqrt{2 \left(1 + \frac{1}{k-1} \right)} = \sqrt{3} \approx 1.732, \quad (32)$$

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

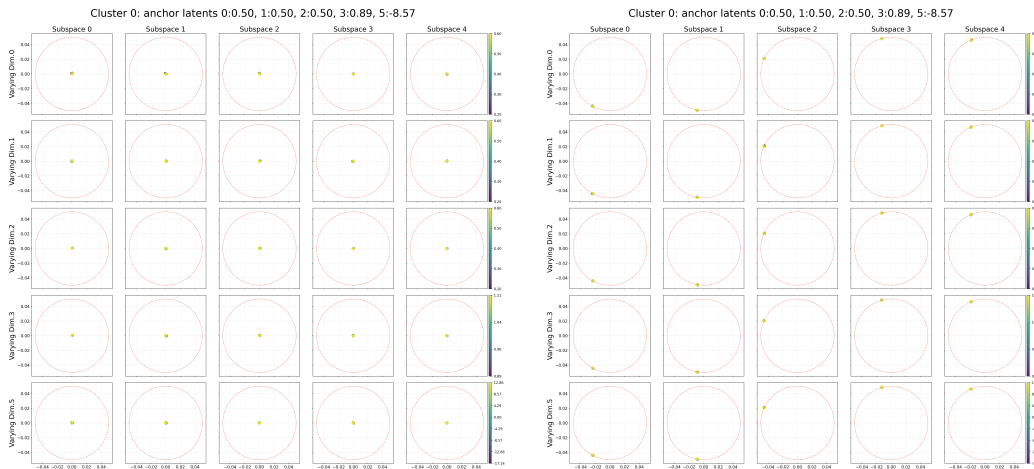


Figure 6: Visualization of Shape3D per aligned rotation subspace. Left: w/o $\mathcal{L}_{\text{radius}}$, all the features collapse to the rotation center. Right: w/o $\mathcal{L}_{\text{theta}}$, all subspaces collapse to a point, i.e., the learned rotation matrix is an identity matrix.

matching the theoretical maximum for three points on a sphere (Tammes, 1930). With fixed (unlearned) rotation centers, we obtain $d_{\min} = 1.719$.

F FULL EXPERIMENT RESULTS

F.1 EXTENDED RESULTS FOR SECTION 4.2

Table 4: DCI using Linear regression on four datasets shown compactly as triples **D/C/I** = Disentanglement/Completeness/Informativeness.

Method/Dataset	MPI3D (D/C/I)	Shape3D (D/C/I)	3DIdent (D/C/I)	3DIEBench (D/C/I)
<i>(Self-Supervised Invariance)</i>				
SimCLR	0.234 / 0.078 / 0.144	0.996 / 0.099 / 0.167	0.286 / 0.084 / 0.313	0.231 / 0.099 / 0.173
<i>(Self-Supervised Equivariance)</i>				
EquiMOD	0.545 / 0.120 / 0.203	0.851 / 0.168 / 0.483	0.232 / 0.078 / 0.309	0.736 / 0.115 / 0.289
CARE	0.236 / 0.111 / 0.169	0.842 / 0.066 / 0.181	0.243 / 0.080 / 0.329	0.174 / 0.096 / 0.309
<i>(Self-Supervised Disentanglement)</i>				
IP-IRM	0.654 / 0.147 / 0.346	0.357 / 0.114 / 0.132	0.202 / 0.070 / 0.431	0.534 / 0.122 / 0.250
Eastwood et al. (2023)	0.950 / 0.700 / 0.414	0.962 / 0.821 / 0.694	0.879 / 0.693 / 0.713	0.948 / 0.501 / 0.276
CD-RED (before post-processing)	0.628 / 0.505 / 1.000	0.630 / 0.494 / 1.000	0.599 / 0.471 / 0.996	0.106 / 0.067 / 0.309
CD-RED (after post-processing)	0.951 / 0.952 / 0.998	0.986 / 0.979 / 0.999	0.980 / 0.980 / 0.996	0.986 / 0.992 / 1.000

Table 5: Per-dimension Informativeness using linear regression on **Shapes3D**. Columns list ground-truth factors

Method	Floor Hue	Wall Hue	Object Hue	Scale	Orientation	Class
SimCLR	0.0000	0.0000	0.0000	0.0000	0.0000	1.0000
EquiMOD	0.3071	0.5003	0.1089	0.0005	0.9789	0.9997
CARE	0.0158	0.0191	0.0185	0.0106	0.0216	0.9986
IP-IRM	0.0279	0.0260	0.0049	0.0259	0.9949	0.9967
Eastwood et al. (2023)	0.5410	0.8655	0.6672	0.6553	0.0003	0.7868
CD-RED (before post-proc.)	0.6860	0.8241	0.9620	0.5743	0.8592	0.9998
CD-RED (after post-proc.)	0.9995	0.9995	0.9994	0.9984	0.9994	1.0000

Table 6: Per-dimension Informativeness using Random forest regression on **Shapes3D**. Columns list ground-truth factors

Method	Floor Hue	Wall Hue	Object Hue	Scale	Orientation	Class
SimCLR	0.8637	0.8719	0.8598	0.8582	0.8581	1.0000
EquiMOD	0.9798	0.9836	0.9497	0.8953	0.9998	1.0000
CARE	0.9200	0.9198	0.9011	0.9093	0.8968	1.0000
IP-IRM	0.9327	0.9415	0.9035	0.9489	1.0000	1.0000
Eastwood et al. (2023)	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
CD-RED (before post-proc.)	1.0000	1.0000	0.9999	1.0000	1.0000	1.0000
CD-RED (after post-proc.)	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000

Table 7: Per-dimension Informativeness using linear regression on **MPI3D**. Columns list the dataset’s ground-truth factors.

Method	Obj. Color	Obj. Size	Cam. Height	Backg. Color	Arm Horiz.	Arm Vert.	Class
SimCLR	0.0005	0.0011	0.0017	0.0008	0.0003	0.0008	1.0000
EquiMOD	0.4191	0.0001	0.0002	0.0008	0.0002	0.0002	0.9998
CARE	0.0323	0.0562	0.0749	0.0118	0.0021	0.0059	0.9999
IP-IRM	0.7589	0.0300	0.0205	0.0015	0.0018	0.0028	0.1049]
Eastwood et al. (2023)	0.5896	0.4471	1.0000	1.0000	0.9996	0.8018	0.0199]
CD-RED (before post-proc.)	0.3758	0.0466	0.2571	0.3304	0.0552	0.8311	0.8388
CD-RED (after post-proc.)	0.9953	0.9983	0.9984	0.9950	0.9977	0.9982	1.0000

Table 8: Per-dimension Informativeness using Random forest regression on **MPI3D**. Columns list the dataset’s ground-truth factors.

Method	Obj. Color	Obj. Size	Cam. Height	Backg. Color	Arm Horiz.	Arm Vert.	Class
SimCLR	0.8575	0.8580	0.8576	0.8562	0.8572	0.8564	1.0000
EquiMOD	0.9911	0.9916	0.9967	0.9962	0.8724	0.8614	0.9999
CARE	0.8795	0.8772	0.8885	0.8597	0.8648	0.8667	1.0000
IP-IRM	0.9881	0.9527	0.9785	0.9768	0.8822	0.8877	0.9994
Eastwood et al. (2023)	1.0000	0.9992	1.0000	1.0000	1.0000	1.0000	0.9994
CD-RED (before post-proc.)	0.9995	0.9999	0.9998	1.0000	0.9997	0.9998	1.0000
CD-RED (after post-proc.)	0.9999	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000

Table 9: Per-dimension Informativeness using Linear Regression on **3DIdent**. Columns follow the latents: Position (X, Y, Z), Rotation (ϕ, θ, ψ), Spotlight Position, and Hue (Object, Ground, Spotlight).

Method	Position			Rotation			Spotlight Pos.	Hue		
	X	Y	Z	ϕ	θ	ψ	Position	Obj.	Ground	Spotlight
SimCLR	0.3728	0.0121	0.2591	0.0005	0.0037	0.0006	0.2361	0.4108	0.1706	0.9747
EquiMOD	0.4786	0.0884	0.4086	0.0008	0.0304	0.0021	0.5336	0.5095	0.2511	0.0996
CARE	0.4726	0.0352	0.3177	0.0007	0.0061	0.0013	0.2237	0.4929	0.1067	0.9599
IP-IRM	0.5563	0.0937	0.6513	0.0013	0.0270	0.0037	0.4954	0.7152	0.2584	0.9309
Eastwood et al. (2023)	0.7366	0.6961	0.7102	0.8021	0.7524	0.8062	0.7258	0.6957	0.5405	0.3730
CD-RED (before post-proc.)	0.9965	0.9971	0.9971	0.9985	0.9982	0.9988	0.9943	0.9805	0.9904	0.9981
CD-RED (after post-proc.)	0.9968	0.9974	0.9974	0.9987	0.9983	0.9989	0.9947	0.9809	0.9908	0.9984

Table 10: Per-dimension Informativeness using Random forest regression on **3DIdent**. Columns follow the latents: Position (X, Y, Z), Rotation (ϕ, θ, ψ), Spotlight Position, and Hue (Object, Ground, Spotlight).

Method	Position			Rotation			Spotlight Pos.	Hue		
	X	Y	Z	ϕ	θ	ψ	Position	Obj.	Ground	Spotlight
SimCLR	0.9599	0.8738	0.9521	0.8533	0.8567	0.8517	0.9457	0.9629	0.9326	0.9993
EquiMOD	0.9584	0.8824	0.9651	0.8558	0.8627	0.8562	0.9535	0.9802	0.9536	0.9956
CARE	0.9554	0.8741	0.9555	0.8543	0.8564	0.8540	0.9419	0.9796	0.9349	0.9994
IP-IRM	0.9546	0.8799	0.9655	0.8554	0.8616	0.8561	0.9541	0.9813	0.9468	0.9923
Eastwood et al. (2023)	1.0000	1.0000	0.9999	0.9997	0.9999	0.9997	0.9999	0.9938	0.9969	0.9666
CD-RED (before post-proc.)	0.9995	0.9996	0.9996	0.9998	0.9997	0.9998	0.9992	0.9977	0.9989	0.9998
CD-RED (after post-proc.)	0.9996	0.9996	0.9996	0.9998	0.9998	0.9998	0.9993	0.9978	0.9989	0.9998

2160

2161

2162

2163

2164 Table 11: Per-dimension Informativeness on **3DIEBench** using Linear Regression. Columns list the
2165 dataset’s ground-truth factors.

2166

2167

Method	Position X	Position Y	Position Z	Rotation ψ	Ground Hue	Class
SimCLR	0.0007	0.0027	0.0057	0.0014	0.0299	0.9997
EquiMOD	0.0008	0.0007	0.0016	0.0006	0.7297	1.0000
CARE	0.0011	0.0099	0.0131	0.0006	0.8314	0.9996
IP-IRM	0.0080	0.1693	0.3599	0.0200	0.0394	0.8489
Eastwood et al. (2023)	0.0011	0.2802	0.2391	0.0173	0.3762	0.7424
CD-RED (before post-proc.)	0.4318	0.0687	0.0600	0.2215	0.0069	0.9993
CD-RED (after post-proc.)	0.9989	0.9988	0.9988	0.9987	0.9994	1.0000

2174

2175

2176

2177

2178

2179

2180

2181

2182 Table 12: Per-dimension Informativeness on **3DIEBench** using Random forest Regression.
2183 Columns list the dataset’s ground-truth factors.

2184

2185

Method	Position X	Position Y	Position Z	Rotation ψ	Ground Hue	Class
SimCLR	0.8564	0.8738	0.8805	0.8635	0.9475	1.0000
EquiMOD	0.8555	0.8730	0.8725	0.8697	0.9977	1.0000
CARE	0.8687	0.9376	0.9495	0.9133	0.9897	1.0000
IP-IRM	0.8696	0.9353	0.9527	0.9118	0.9046	1.0000
Eastwood et al. (2023)	0.8654	1.0000	0.9999	0.8958	0.9502	1.0000
CD-RED (before post-proc.)	0.9997	0.9996	0.9997	0.9997	0.9998	1.0000
CD-RED (after post-proc.)	0.9998	0.9998	0.9998	0.9998	0.9999	1.0000

2192

2193

2194

2195

2196

2197

2198

2199

2200

2201

2202

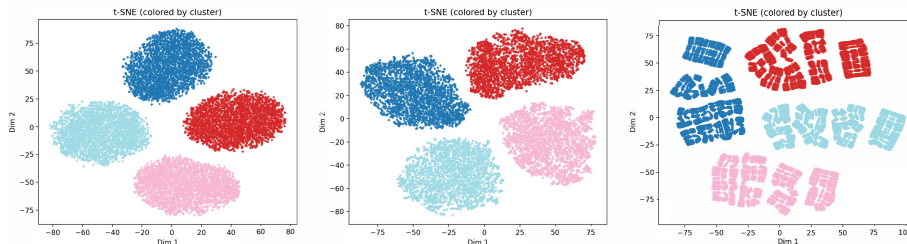
2203

2204

2205

2206

2207



2208

2209 Figure 7: t-SNE visualization of learned features on **Shapes3D**. **Left**: Simclr; **Middle**: CARE;
2210 **Right**: Ours.

2211

2212

2213

2214
 2215
 2216
 2217
 2218
 2219
 2220
 2221
 2222
 2223
 2224
 2225
 2226
 2227
 2228
 2229
 2230
 2231
 2232
 2233
 2234
 2235
 2236
 2237
 2238
 2239
 2240
 2241
 2242
 2243
 2244
 2245
 2246
 2247
 2248
 2249
 2250
 2251
 2252
 2253
 2254
 2255
 2256
 2257
 2258
 2259
 2260
 2261
 2262
 2263
 2264
 2265
 2266
 2267

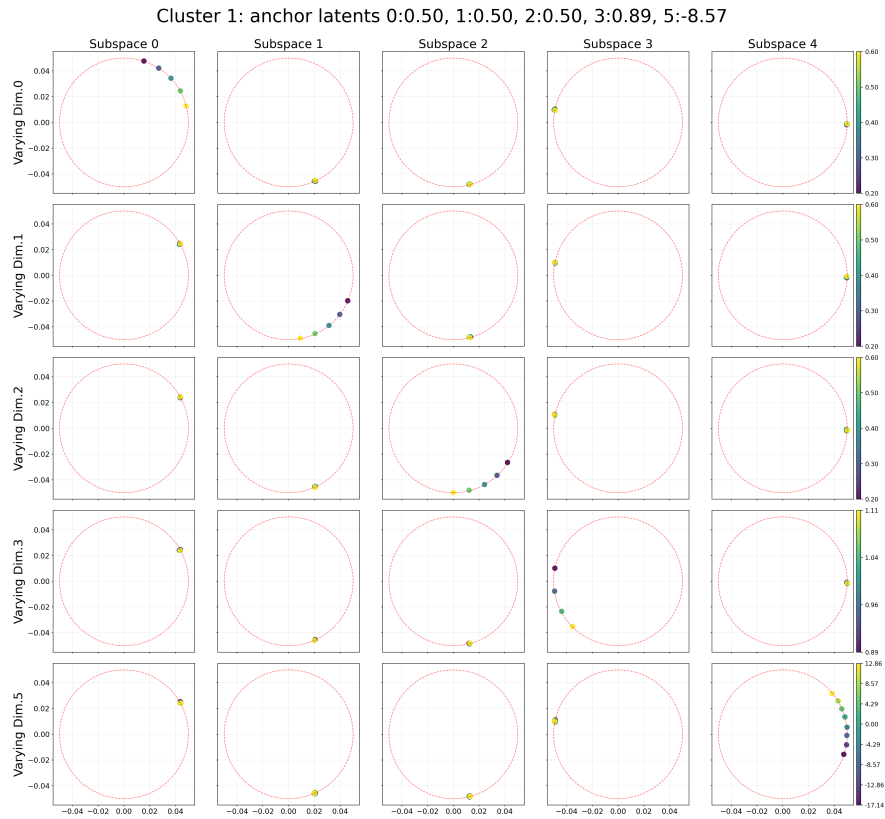


Figure 8: Randomly pick an anchor in one cluster and visualize each latent dimension varying on **Shapes3D**. Dim 0: Floor Hue, Dim 1: Wall Hue, Dim 2: Object Hue Dim 3: Scale Dim 5: Orientation

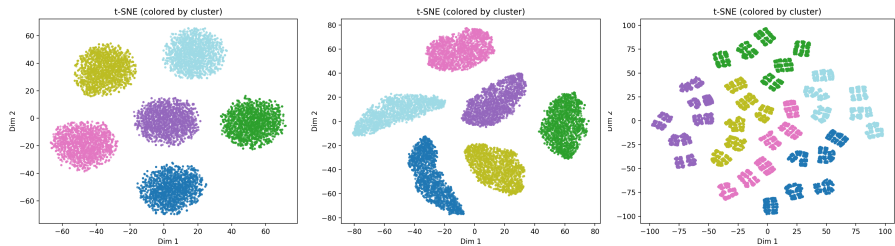


Figure 9: t-SNE visualization of learned features on **MPI3D**. **Left:** Simclr; **Middle:** CARE; **Right:** Ours.

2268
 2269
 2270
 2271
 2272
 2273
 2274
 2275
 2276
 2277
 2278
 2279
 2280
 2281
 2282
 2283
 2284
 2285
 2286
 2287
 2288
 2289
 2290
 2291
 2292
 2293
 2294
 2295
 2296
 2297
 2298
 2299
 2300
 2301
 2302
 2303
 2304
 2305
 2306
 2307
 2308
 2309
 2310
 2311
 2312
 2313
 2314
 2315
 2316
 2317
 2318
 2319
 2320
 2321

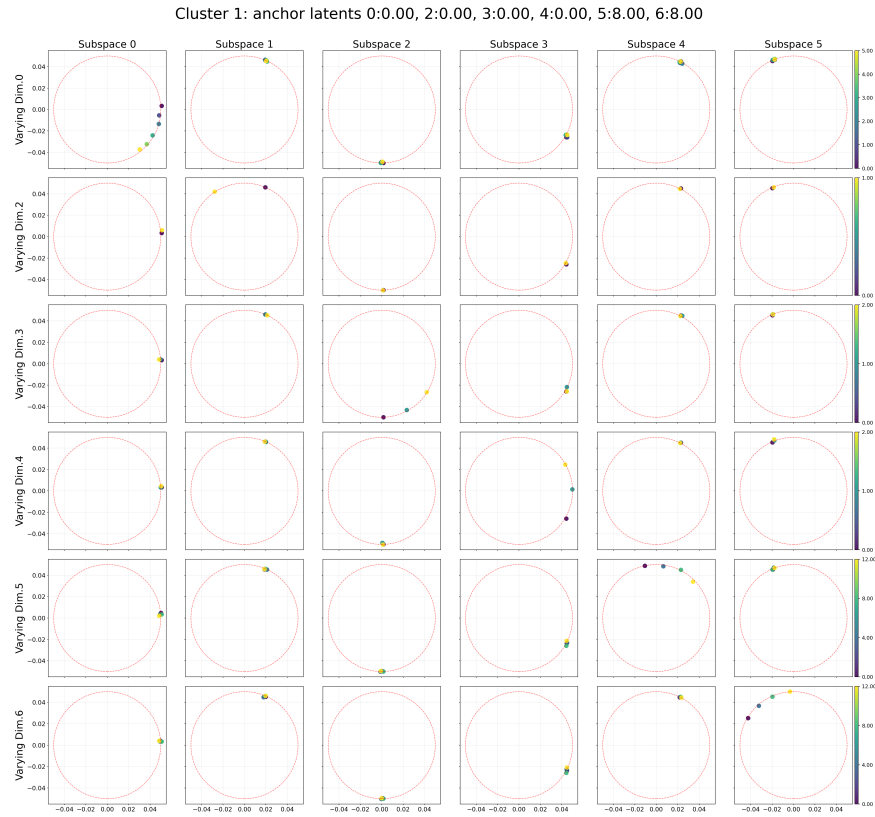


Figure 10: Randomly pick an anchor in one cluster and visualize each latent dimension varying on **MPI3D**. Dim 0: Floor Hue, Dim 1: Wall Hue, Dim 2: Object Hue, Dim 3: Scale, Dim 5: Orientation

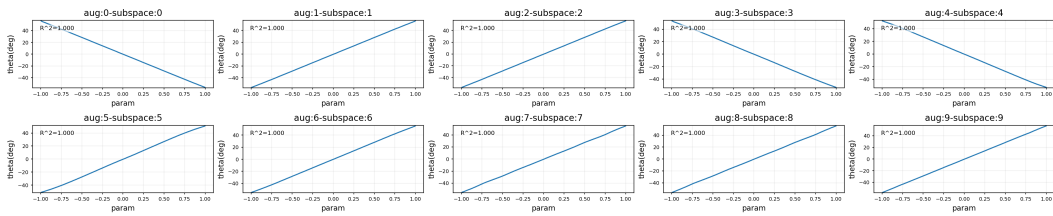


Figure 11: Transformation network output for each augmentation in **3DIdent**. Randomly pick 1000 points within the augmentation range.

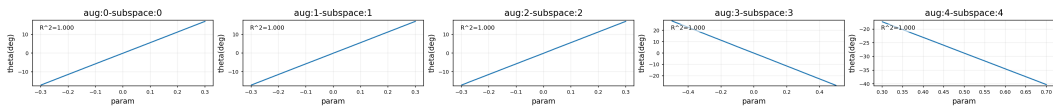


Figure 12: Transformation network output for each augmentation **3DIEBench**. Randomly pick 1000 points within the augmentation range.

F.2 EXTENDED RESULTS FOR SECTION 4.3

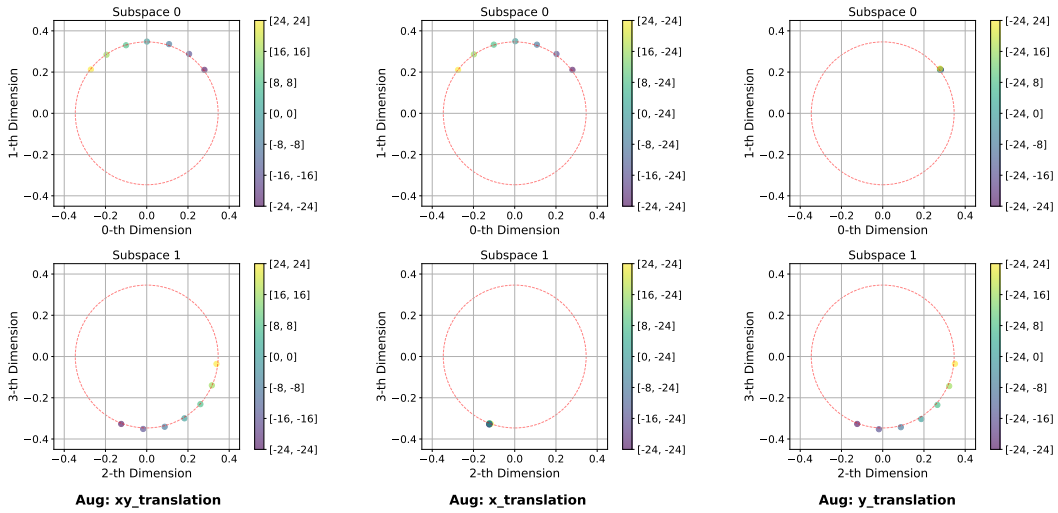


Figure 13: Experiment with x-translation, y-translation, and xy-translation. It has effectively disentangled the x-translation and y-translation, while the xy-translation learns to perfectly decompose and align its caused variations to the two subspaces assigned to the x-translation and y-translation.

2376
 2377
 2378
 2379
 2380
 2381
 2382
 2383
 2384
 2385
 2386
 2387
 2388
 2389
 2390
 2391
 2392
 2393
 2394
 2395
 2396
 2397
 2398
 2399
 2400
 2401
 2402
 2403
 2404
 2405
 2406
 2407
 2408
 2409
 2410
 2411
 2412
 2413
 2414
 2415
 2416
 2417
 2418
 2419
 2420
 2421
 2422
 2423
 2424
 2425
 2426
 2427
 2428
 2429

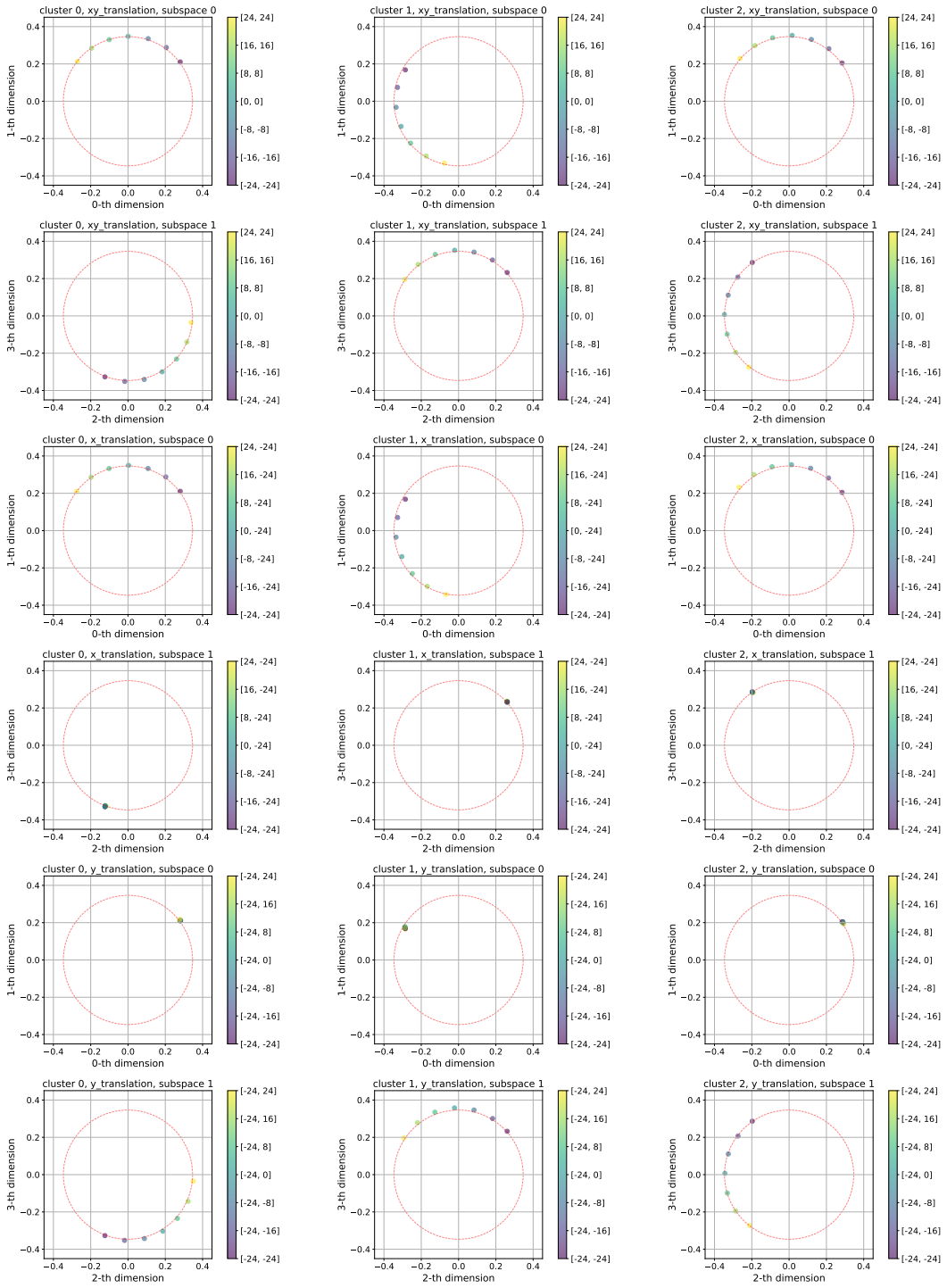


Figure 14: Full visualization: x-translation, y-translation, xy-translation.

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

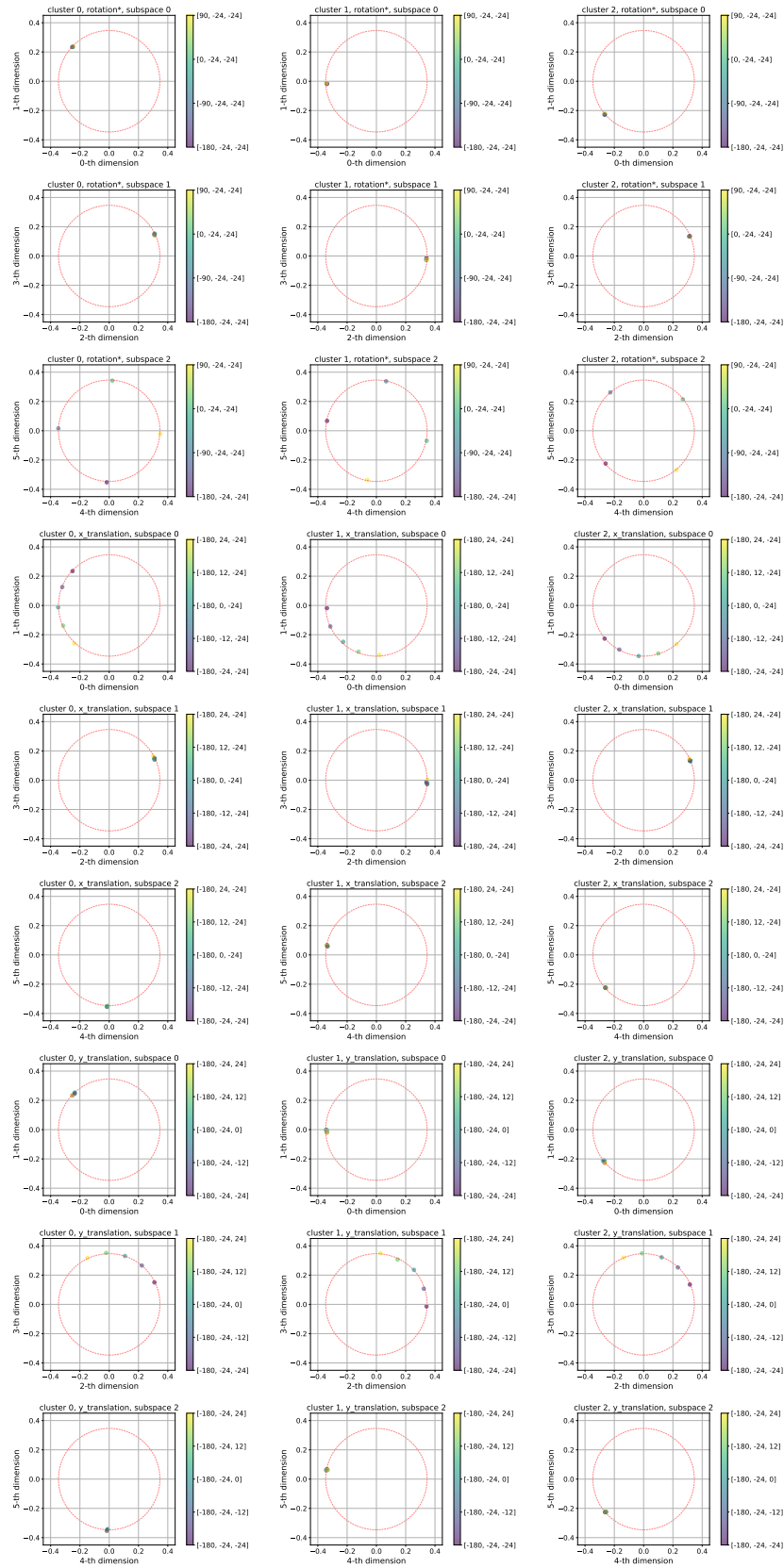


Figure 15: Full visualization: x-translation, y-translation, image-centered rotation.