

# Equivariant Representation Learning in the Presence of Stabilizers

Luis Armando Pérez Rey (✉) <sup>\*1,2,3</sup>, Giovanni Luca Marchetti <sup>\*4</sup>,  
Danica Kragic<sup>4</sup>, Dmitri Jarnikov<sup>1,3</sup>, and Mike Holenderski<sup>1</sup>

<sup>1</sup> Eindhoven University of Technology, Eindhoven, The Netherlands

<sup>2</sup> Eindhoven Artificial Intelligence Systems Institute, Eindhoven, The Netherlands

<sup>3</sup> Prosus, Amsterdam, The Netherlands

<sup>4</sup> KTH Royal Institute of Technology, Stockholm, Sweden

**Abstract.** We introduce Equivariant Isomorphic Networks (EquIN) – a method for learning representations that are equivariant with respect to general group actions over data. Differently from existing equivariant representation learners, EquIN is suitable for group actions that are not free, i.e., that stabilize data via nontrivial symmetries. EquIN is theoretically grounded in the orbit-stabilizer theorem from group theory. This guarantees that an ideal learner infers isomorphic representations while trained on equivariance alone and thus fully extracts the geometric structure of data. We provide an empirical investigation on image datasets with rotational symmetries and show that taking stabilizers into account improves the quality of the representations.

**Keywords:** Representation Learning · Equivariance · Lie Groups

## 1 Introduction

Incorporating data symmetries into deep neural representations defines a fundamental challenge and has been addressed in several recent works [1, 6, 14, 26, 28]. The overall aim is to design representations that preserve symmetries and operate coherently with respect to them – a functional property known as *equivariance*. This is because the preservation of symmetries leads to the extraction of geometric and semantic structures in data, which can be exploited for data efficiency and generalization [2]. As an example, the problem of *disentangling* semantic factors of variation in data has been rephrased in terms of equivariant representations [3, 12]. As disentanglement is known to be unfeasible with no inductive biases or supervision [18], symmetries of data arise as a geometric structure that can provide weak supervision and thus be leveraged in order to disentangle semantic factors.

The majority of models from the literature rely on the assumption that the group of symmetries acts *freely* on data [21] i.e., that no datapoint is stabilized by nontrivial symmetries. This avoids the need to model *stabilizers* of datapoints,

---

\*Equal Contribution

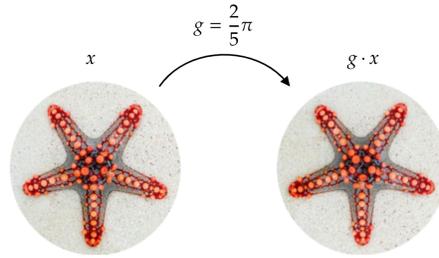


Fig. 1: An example of an action on data that is not free. The datapoint  $x$  is stabilized by the symmetry  $g \in G$ .

which are unknown subgroups of the given symmetry group. However, non-free group actions arise in several practical scenarios. This happens, for example, when considering images of objects acted upon by the rotation group via a change of orientation. Such objects may be symmetrical, resulting in rotations leaving the image almost identical and consequently ambiguous in its orientation, see Figure 1. Discerning the correct orientations of an object is important for applications such as pose estimation [20] and reinforcement learning [10]. This motivates the need to design equivariant representation learning frameworks that are capable of modeling stabilizers and therefore suit non-free group actions.

In this work, we propose a method for learning equivariant representation for general and potentially non-free group actions. Based on the Orbit-Stabilizer Theorem from group theory, we design a model that outputs subsets of the group, which represent the stabilizer subgroup up to a symmetry – a group theoretical construction known as *coset*. The representation learner optimizes an equivariance loss relying on supervision from symmetries alone. This means that we train our model on a dataset consisting of relative symmetries between pairs of datapoints, avoiding the need to know the whole group action over data a priori. From a theoretical perspective, the above-mentioned results from group theory guarantee that an ideal learner infers representations that are isomorphic to the original dataset. This implies that our representations completely preserve the symmetry structure while preventing any loss of information. We name our framework Equivariant Isomorphic Networks – EquIN for short. In summary, our contributions include:

- A novel equivariant representation learning framework suitable for non-free group actions.
- A discussion grounded on group theory with theoretical guarantees for isomorphism representations.
- An empirical investigation with comparisons to competing equivariant representation learners on image datasets.

We provide Python code implementing our framework together with all the experiments at the following repository: [luis-armando-perez-rey/non-free](https://github.com/luis-armando-perez-rey/non-free).

## 2 Related Work

In this section, we first briefly survey representation learning methods from the literature leveraging on equivariance. We then draw connections between equivariant representations and world models from reinforcement learning and discuss the role of equivariance in terms of disentangling semantic factors of data.

**Equivariant Representation Learning.** Several works in the literature have proposed and studied representation learning models that are equivariant with respect to a group of data symmetries. These models are typically trained via a loss encouraging equivariance on a dataset of relative symmetries between datapoints. What distinguishes the models is the choice of the latent space and of the group action over the latter. Euclidean latent spaces with linear or affine actions have been explored in [9, 26, 29]. However, the intrinsic data manifold is non-Euclidean in general, leading to representations that are non-isomorphic and that do not preserve the geometric structure of the data. To amend this, a number of works have proposed to design latent spaces that are isomorphic to disjoint copies of the symmetry group [8, 15, 21, 28]. When the group action is free, this leads to isomorphic representations and thus completely recovers the geometric structure of the data [21]. However, the proposed latent spaces are unsuitable for non-free actions. Since they do not admit stabilizers, no equivariant map exists, and the model is thus unable to learn a suitable representation. In the present work, we extend this line of research by designing a latent space that enables learning equivariant representations in the presence of stabilizers. Our model implicitly represents stabilizer subgroups and leads to isomorphic representations for arbitrary group actions.

**Latent World Models.** Analogously to group actions, Markov Decision Processes (MDPs) from reinforcement learning and control theory involve a, possibly stochastic, interaction with an environment. This draws connections between MDPs and symmetries since the latter can be thought of as transformations and, thus, as a form of interaction. The core difference is that in an MDP, no algebraic structure, such as a group composition, is assumed on the set of interactions. In the context of MDPs, a representation that is equivariant with respect to the agent’s actions is referred to as latent *World Model* [10, 16, 24] or *Markov Decision Process Homomorphism* (MDPH) [25]. In an MDPH the latent action is learned together with the representation by an additional model operating on the latent space. Although this makes MDPHs more general than group-equivariant models, the resulting representation is unstructured and uninterpretable. The additional assumptions of equivariant representations translate instead into the preservation of the geometric structure of data.

**Disentanglement.** As outlined in [2], a desirable property for representations is disentanglement, i.e., the ability to decompose in the representations the semantic factors of variations that explain the data. Although a number of methods have been proposed for this purpose [4, 13], it has been shown that disentanglement is mathematically unachievable in an unbiased and unsupervised way [18]. As an alternative, the notion has been rephrased in terms of symmetry and equivariance [12]. It follows that isomorphic equivariant representations are

guaranteed to be disentangled in this sense [21, 28]. Since we aim for general equivariant representations that are isomorphic, our proposed method achieves disentanglement as a by-product.

### 3 Group Theory Background

We review the fundamental group theory concepts necessary to formalize our representation learning framework. For a complete treatment, we refer to [27].

**Definition 1.** *A group is a set  $G$  equipped with a composition map  $G \times G \rightarrow G$  denoted by  $(g, h) \mapsto gh$ , an inversion map  $G \rightarrow G$  denoted by  $g \mapsto g^{-1}$ , and a distinguished identity element  $1 \in G$  such that for all  $g, h, k \in G$ :*

$$\begin{array}{lll} \text{Associativity} & \text{Inversion} & \text{Identity} \\ g(hk) = (gh)k & g^{-1}g = gg^{-1} = 1 & g1 = 1g = g \end{array}$$

Elements of a group represent abstract symmetries. Spaces with a group of symmetries  $G$  are said to be acted upon by  $G$  in the following sense.

**Definition 2.** *An action by a group  $G$  on a set  $\mathcal{X}$  is a map  $G \times \mathcal{X} \rightarrow \mathcal{X}$  denoted by  $(g, x) \mapsto g \cdot x$ , satisfying for all  $g, h \in G$ ,  $x \in \mathcal{X}$ :*

$$\begin{array}{ll} \text{Associativity} & \text{Identity} \\ g \cdot (h \cdot x) = (gh) \cdot x & 1 \cdot x = x \end{array}$$

Suppose that  $G$  acts on a set  $\mathcal{X}$ . The action defines a set of *orbits*  $\mathcal{X}/G$  given by the equivalence classes of the relation  $x \sim y$  iff  $y = g \cdot x$  for some  $g \in G$ . For each  $x \in \mathcal{X}$ , the *stabilizer* subgroup is defined as

$$G_x = \{g \in G \mid g \cdot x = x\}. \quad (1)$$

Stabilizers of elements in the same orbit are conjugate, meaning that for each  $x, y$  belonging to the same orbit  $O$  there exists  $h \in G$  such that  $G_y = hG_xh^{-1}$ . By abuse of notation, we refer to the conjugacy class  $G_O$  of stabilizers for  $O \in \mathcal{X}/G$ . The action is said to be *free* if all the stabilizers are trivial, i.e.,  $G_O = \{1\}$  for every  $O$ .

We now recall the central notion for our representation learning framework.

**Definition 3.** *A map  $\varphi : \mathcal{X} \rightarrow \mathcal{Z}$  between sets acted upon by  $G$  is equivariant if  $\varphi(g \cdot x) = g \cdot \varphi(x)$  for every  $x \in \mathcal{X}$  and  $g \in G$ . An equivariant bijection is referred to as isomorphism.*

Intuitively, an equivariant map between  $\mathcal{X}$  and  $\mathcal{Z}$  preserves their corresponding symmetries. The following is the fundamental result on group actions [27].

**Theorem 1 (Orbit-Stabilizer).** *The following holds:*

- Each orbit  $O$  is isomorphic to the set of (left) cosets  $G/G_O = \{gG_O \mid g \in G\}$ . In other words, there is an isomorphism:

$$\mathcal{X} \simeq \coprod_{O \in \mathcal{X}/G} G/G_O \subseteq 2^G \times \mathcal{X}/G \quad (2)$$

where  $2^G$  denotes the power-set of  $G$  on which  $G$  acts by left multiplication i.e.,  $g \cdot A = \{ga \mid a \in A\}$ .

- Any equivariant map

$$\varphi: \mathcal{X} \rightarrow \coprod_{O \in \mathcal{X}/G} G/G_O \quad (3)$$

that induces a bijection on orbits is an isomorphism.

Theorem 1 describes arbitrary group actions completely and asserts that orbit-preserving equivariant maps are isomorphisms. Our central idea is to leverage on this in order to design a representation learner that is guaranteed to be isomorphic when trained on equivariance alone.

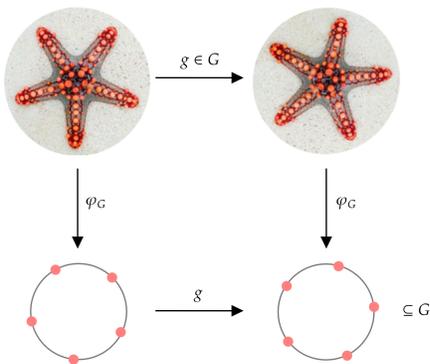


Fig. 2: An illustration of EquIN encoding data equivariantly as subsets of the symmetry group  $G$ . This results in representations that are suitable even when the action by  $G$  on data is not free.

## 4 Equivariant Isomorphic Networks (EquIN)

Our goal is to design an equivariant representation learner based on Theorem 1. We aim to train a model

$$\varphi: \mathcal{X} \rightarrow \mathcal{Z} \quad (4)$$

with a latent space  $\mathcal{Z}$  on a loss encouraging equivariance. The ideal choice for  $\mathcal{Z}$  is given by  $\coprod_{O \in \mathcal{X}/G} G/G_O$  since the latter is isomorphic to  $\mathcal{X}$  (Theorem 1).

In other words,  $\varphi$  ideally outputs cosets of stabilizers of the input datapoints. However, while we assume that  $G$  is known a priori, its action on  $\mathcal{X}$  is not and has to be inferred from data. Since the stabilizers depend on the group action, they are unknown a priori as well. In order to circumvent the modeling of stabilizers and their cosets, we rely on the following simple result:

**Proposition 1.** *Let  $\varphi : \mathcal{X} \rightarrow 2^G$  be an equivariant map. Then for each  $x \in \mathcal{X}$  belonging to an orbit  $O$ ,  $\varphi(x)$  contains a coset of (a conjugate of)  $G_O$ .*

*Proof.* Pick  $x \in \mathcal{X}$ . Then for every  $g \in G_x$  it holds that  $\varphi(x) = \varphi(g \cdot x) = g \cdot \varphi(x)$ . In other words,  $G_x h = h h^{-1} G_x h \subseteq \varphi(x)$  for each  $h \in \varphi(x)$ . Since  $h^{-1} G_x h$  is conjugate to  $G_x$  the thesis follows.

Proposition 1 enables  $\varphi$  to output arbitrary subsets of  $G$  instead of cosets of stabilizers. As long as those subsets are *minimal* w.r.t. to inclusion, they will coincide with the desired cosets.

Based on this, we define the latent space of EquIN as  $\mathcal{Z} = \mathcal{Z}_G \times \mathcal{Z}_O$  and implement the map  $\varphi$  as a pair of neural networks  $\varphi_G : \mathcal{X} \rightarrow \mathcal{Z}_G$  and  $\varphi_O : \mathcal{X} \rightarrow \mathcal{Z}_O$ . The component  $\mathcal{Z}_G$  represents cosets of stabilizers while  $\mathcal{Z}_O$  represents orbits. Since the output space of a neural network is finite-dimensional, we assume that the stabilizers of the action are finite. The model  $\varphi_G$  then outputs  $N$  elements

$$\varphi_G(x) = \{\varphi_G^1(x), \dots, \varphi_G^N(x)\} \subseteq G \quad (5)$$

where  $\varphi_G^i(x) \in G$  for all  $i$ . The hyperparameter  $N$  should be ideally chosen larger than the cardinality of the stabilizers. On the other hand, the output of  $\varphi_O$  consists of a vector of arbitrary dimensionality. The only requirement is that the output space of  $\varphi_O$  should have enough capacity to contain the space of orbits  $\mathcal{X}/G$ .

#### 4.1 Parametrizing $G$ via the Exponential Map

The output space of usual machine learning models such as deep neural networks is Euclidean. However,  $\varphi_G$  needs to output elements of the group  $G$  (see Equation 5), which may be non-Euclidean as in the case of  $G = \text{SO}(n)$ . Therefore, in order to implement  $\varphi_G$ , it is necessary to parametrize  $G$ . To this end, we assume that  $G$  is a differentiable manifold, with differentiable composition and inversion maps, i.e., that  $G$  is a *Lie group*. One can then define the *Lie algebra*  $\mathfrak{g}$  of  $G$  as the tangent space to  $G$  at 1.

We propose to rely on the *exponential map*  $\mathfrak{g} \rightarrow G$ , denoted by  $v \mapsto e^v$ , to parametrize  $G$ . This means that  $\varphi_G$  first outputs  $N$  elements  $\varphi_G(x) = \{v^1, \dots, v^N\} \subseteq \mathfrak{g}$  that get subsequently mapped into  $G$  as  $\{e^{v^1}, \dots, e^{v^N}\}$ . Although the exponential map can be defined for general Lie groups by solving an appropriate ordinary differential equation, we focus on the case  $G \subseteq \text{GL}(n)$ . The

Lie algebra  $\mathfrak{g}$  is then contained in the space of  $n \times n$  matrices and the exponential map amounts to the matrix Taylor expansion

$$e^v = \sum_{k \geq 0} \frac{v^k}{k!} \quad (6)$$

where  $v^k$  denotes the power of  $v$  as a matrix. For specific groups, the latter can be simplified via simple closed formulas. For example, the exponential map of  $\mathbb{R}^n$  is the identity while for  $\text{SO}(3)$  it can be efficiently computed via the Rodrigues' formula [17].

## 4.2 Training Objective

As mentioned, our dataset  $\mathcal{D}$  consists of samples from the unknown group action. This means that datapoints are triplets  $(x, g, y) \in \mathcal{X} \times G \times \mathcal{X}$  with  $y = g \cdot x$ . Given a datapoint  $(x, g, y) \in \mathcal{D}$  the learner  $\varphi_G$  optimizes the equivariance loss over its parameters:

$$\mathcal{L}_G(x, g, y) = d(g \cdot \varphi_G(x), \varphi_G(y)) \quad (7)$$

where  $d$  is a semi-metric for sets. We opt for the asymmetric *Chamfer distance*

$$d(A, B) = \frac{1}{|A|} \sum_{a \in A} \min_{b \in B} d_G(a, b) \quad (8)$$

because of its differentiability properties. Any other differentiable distance between sets of points can be deployed as an alternative. Here  $d_G$  is a metric on  $G$  and is typically set as the squared Euclidean for  $G = \mathbb{R}^n$  and as the squared Frobenius for  $G = \text{SO}(n)$ . As previously discussed, we wish  $\varphi_G(x)$ , when seen as a set, to be minimal in cardinality. To this end, we add the following regularization term measuring the discrete entropy:

$$\tilde{\mathcal{L}}_G(x) = \frac{\lambda}{N^2} \sum_{1 \leq i, j \leq N} d_G(\varphi_G^i(x), \varphi_G^j(x)) \quad (9)$$

where  $\lambda$  is a weighting hyperparameter. On the other hand, since orbits are invariant to the group action  $\varphi_O$  optimizes a *contrastive loss*. We opt for the popular InfoNCE loss from the literature [5]:

$$\mathcal{L}_O(x, y) = d_O(\varphi_O(x), \varphi_O(y)) + \log \mathbb{E}_{x'} \left[ e^{-d_O(\varphi_O(x'), \varphi_O(x))} \right] \quad (10)$$

where  $x'$  is marginalized from  $\mathcal{D}$ . As customary for the InfoNCE loss, we normalize the output of  $\varphi_O$  and set  $d_O(a, b) = -\cos(\angle ab) = -a \cdot b$ . The second summand of  $\mathcal{L}_O$  encourages injectivity of  $\varphi_O$  and as such prevents orbits from overlapping in the representation.

The Orbit-Stabilizer Theorem (Theorem 1) guarantees that if EquIN is implemented with ideal learners  $\varphi_G, \varphi_O$  then it infers isomorphic representations in

the following sense. If the  $\mathcal{L}_G(x, g, y)$  and the first summand of  $\mathcal{L}_O(x, y)$  vanish for every  $(x, g, y)$  then  $\varphi$  is equivariant. If moreover the regularizations,  $\tilde{\mathcal{L}}_G$  and the second summand of  $\mathcal{L}_O$ , are at a minimum then  $\varphi_G(x)$  coincides with a coset of  $G_O$  for every  $x \in O$  (Proposition 1) and  $\varphi_O$  is injective. The second claim of Theorem 1 implies then that the representation is isomorphic on its image, as desired.

## 5 Experiments

We empirically investigate EquIN on image data acted upon by a variety Lie groups. Our aim is to show both qualitatively and quantitatively that EquIN reliably infers isomorphic equivariant representations for non-free group actions.

We implement the neural networks  $\varphi_G$  and  $\varphi_O$  as a ResNet18 [11]. For a datapoint  $x \in \mathcal{X}$ , the network implements multiple heads to produce embeddings  $\{\varphi_G^1(x), \dots, \varphi_G^N(x)\} \subseteq G$ . The output dimension of  $\varphi_O$  is set to 3. We train the model for 50 epochs using the AdamW optimizer [19] with a learning rate of  $10^{-4}$  and batches of 16 triplets  $(x, g, y) \in \mathcal{D}$ .

### 5.1 Datasets

We consider the following datasets consisting of  $64 \times 64$  images subject to non-free group actions. Samples from these datasets are shown in Figure 4.

- ROTATING ARROWS: images of radial configurations of  $\nu \in \{1, 2, 3, 4, 5\}$  arrows rotated by  $G = \text{SO}(2)$ . The number of arrows  $\nu$  determines the orbit and the corresponding stabilizer is (isomorphic to) the cyclic group  $C_\nu$  of cardinality  $\nu$ . The dataset contains 2500 triplets  $(x, g, y)$  per orbit.
- COLORED ARROWS: images similar to ROTATING ARROWS but with the arrows of five different colors. This extra factor produces additional orbits with the same stabilizer subgroups. The number of orbits is therefore 25. The dataset contains 2000 triplets per orbit.
- DOUBLE ARROWS: images of two radial configurations of 2, 3 and 3, 5 arrows respectively rotated by the torus  $G = \text{SO}(2) \times \text{SO}(2)$ . The action produces two orbits with stabilizers given by products of cyclic groups:  $C_2 \times C_3$  and  $C_3 \times C_5$  respectively. The dataset contains 2000 triplets per orbit.
- MODELNET: images of monochromatic objects from ModelNet40 [30] rotated by  $G = \text{SO}(2)$  along an axis. We consider five objects: an airplane, a chair, a lamp, a bathtub and a stool. Each object corresponds to a single orbit. The lamp, the stool and the chair have the cyclic group  $C_4$  as stabilizer while the action over the airplane and the bathtub is free. The dataset contains 2500 triplets per orbit.
- SOLIDS: images of a monochromatic tetrahedron, cube and icosahedron [22] rotated by  $G = \text{SO}(3)$ . Each solid defines an orbit, and the stabilizers of the tetrahedron, the cube, and the icosahedron are subgroups of order 12, 24 and 60 respectively. The dataset contains 7500 triplets per orbit.

### 5.2 Comparisons

We compare EquIN with the following two equivariant representation learning models.

- *Baseline*: a model corresponding to EquIN with  $N = 1$  where  $\varphi_G$  outputs a single element of  $G$ . The latent space is  $\mathcal{Z} = G \times \mathcal{Z}_O$ , on which  $G$  acts freely. We deploy this as the baseline since it has been proposed with minor variations in a number of previous works [3, 21, 23, 28] assuming free group actions.
- *Equivariant Neural Renderer (ENR)*: a model from [7] implementing a tensorial latent space  $\mathcal{Z} = \mathbb{R}^{S^3}$ , thought as a scalar signal space on a  $S \times S \times S$  grid in  $\mathbb{R}^3$ . The group  $\text{SO}(3)$  act *approximately* on  $\mathcal{Z}$  by rotating the grid and interpolating the obtained values. The model is trained jointly with a decoder  $\psi : \mathcal{Z} \rightarrow \mathcal{X}$  and optimizes a variation of the equivariance loss that incorporates reconstruction:  $\mathbb{E}_{x, g, y = g \cdot x} [d_{\mathcal{X}}(y, \psi(g \cdot \varphi(x)))]$  where  $d_{\mathcal{X}}$  is the binary cross-entropy for normalized images. Although the action on  $\mathcal{Z}$  is free, the latent discretization and consequent interpolation make the model only approximately equivariant. Similarly to EquIN, we implement ENR as ResNet18. As suggested in the original work [7] we deploy 3D convolutional layers around the latent and set to zero the latent dimensions outside a ball. We set  $S = 8$  with 160 non-zero latent dimensions since this value is comparable to the latent dimensionality of EquIN, between 7 and 250 dimensions depending on  $N$ , making the comparison fair. Note that ENR is inapplicable to DOUBLE ARROWS since its symmetry group is not naturally embedded into  $\text{SO}(3)$ .

### 5.3 Quantitative Results

In order to quantitatively compare the models, we rely on the following evaluation metrics computed on a test dataset  $\mathcal{D}_{\text{test}}$  consisting of 10% of the corresponding training data:

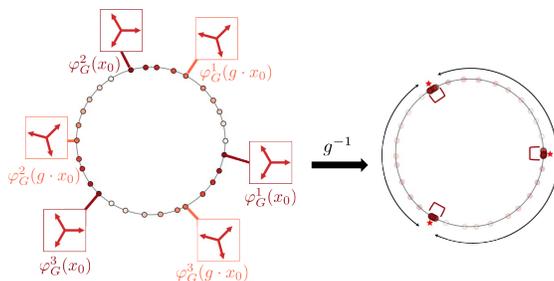


Fig. 3: Diagram explaining the estimation of the disentanglement metric for EquIN. This example assumes that  $G = \text{SO}(2)$  and that  $A$  is the identity.

- *Hit-Rate*: a standard score comparing equivariant representations with different latent space geometries [16]. Given a test triple  $(x, g, y = g \cdot x) \in \mathcal{D}_{\text{test}}$ , we say that ‘ $x$  hits  $y$ ’ if  $\varphi(y)$  is the nearest neighbor in  $\mathcal{Z}$  of  $g \cdot \varphi(x)$  among a random batch of encodings  $\{\varphi(x)\}_{x \in \mathcal{B}}$  with  $|\mathcal{B}| = 20$ . The hit-rate is then defined as the number of times  $x$  hits  $y$  divided by the test set size. For each model, the nearest neighbor is computed with respect to the same latent metric  $d$  as the one used for training. Higher values of the metric are better.
- *Disentanglement*: an evaluation metric proposed in [28] to measure disentanglement according to the symmetry-based definition of [12]. This metric is designed for groups in the form  $G = \text{SO}(2)^T$  and therefore is inapplicable to the SOLIDS dataset. Per orbit, the test set is organized into datapoints of the form  $y = g \cdot x_0$  where  $x_0$  is an arbitrary point in the given orbit. In order to compute the metric, the test dataset is encoded into  $\mathcal{Z}$  via the given representation and then projected to  $\mathbb{R}^{2T}$  via principal component analysis. Then for each independent copy of  $\text{SO}(2) \subseteq G$ , a group action on the corresponding copy of  $\mathbb{R}^2$  is inferred by fitting parameters via a grid search. Finally, the metric computes the average dispersion of the transformed embeddings as the variance of  $g^{-1} \cdot A\varphi_G(y)$ . For EquIN, we propose a modified version accounting for the fact that  $\varphi_G$  produces multiple points in  $G$  using the Chamfer distance  $d$  and averaging the dispersion with respect to each transformed embedding, see Figure 3. The formula for computing the metric is given by:

$$\mathbb{E}_{y,y'}[d(h^{-1} \cdot A\varphi_G(y'), g^{-1} \cdot A\varphi_G(y))] \quad (11)$$

where  $y = g \cdot x_0$  and  $y' = h \cdot x_0$ . Lower values of the metric are better.

The results are summarized in Table 1. EquIN achieves significantly better scores than the baseline. The latter is unable to model the stabilizers in its latent space, leading to representations of poor quality and loss of information. ENR is instead competitive with EquIN. Its latent space suits non-free group actions since stabilizers can be modelled as signals over the latent three-dimensional grid. ENR achieves similar values of hit-rate compared to EquIN. The latter generally outperforms ENR, especially on the MODELNET dataset, while is outperformed on ROTATING ARROWS. According to the disentanglement metric, EquIN achieves significantly lower scores than ENR. This is probably due to the fact the latent group action in ENR is approximate, making the model unable to infer representations that are equivariant at a granular scale.

#### 5.4 Qualitative Results

We provide a number of visualizations as a qualitative evaluation of EquIN. Figure 4 illustrates the output of  $\varphi_G$  on the various datasets. As can be seen, EquIN correctly infers the stabilizers i.e., the cyclic subgroups of  $\text{SO}(2)$  and the subgroup of  $\text{SO}(3)$  of order 12. When  $N$  is larger than the ground-truth cardinalities of stabilizers, the points  $\varphi_G^i$  are overlapped and collapse to the

Table 1: Mean and standard deviation of the metrics across five repetitions. The number juxtaposed to the name of EquIN indicates the cardinality  $N$  of the output of  $\varphi_G$ .

Dataset	Model	Disentanglement ( $\downarrow$ )	Hit-Rate ( $\uparrow$ )
ROTATING ARROWS	Baseline	$1.582_{\pm 0.013}$	$0.368_{\pm 0.004}$
	EquIN5	<b><math>0.009_{\pm 0.005}</math></b>	$0.880_{\pm 0.021}$
	EquIN10	$0.092_{\pm 0.063}$	$0.857_{\pm 0.050}$
	ENR	$0.077_{\pm 0.028}$	<b><math>0.918_{\pm 0.009}</math></b>
COLORED ARROWS	Baseline	$1.574_{\pm 0.007}$	$0.430_{\pm 0.004}$
	EquIN5	$0.021_{\pm 0.015}$	$0.930_{\pm 0.055}$
	EquIN10	<b><math>0.001_{\pm 0.001}</math></b>	<b><math>0.976_{\pm 0.005}</math></b>
DOUBLE ARROWS	Baseline	$1.926_{\pm 0.019}$	$0.023_{\pm 0.004}$
	EquIN6	$0.028_{\pm 0.006}$	$0.512_{\pm 0.011}$
	EquIN15	$0.004_{\pm 0.001}$	$0.820_{\pm 0.104}$
MODELNET	EquIN20	<b><math>0.002_{\pm 0.001}</math></b>	<b><math>0.934_{\pm 0.020}</math></b>
	Baseline	$1.003_{\pm 0.228}$	$0.538_{\pm 0.086}$
	EquIN4	$0.012_{\pm 0.022}$	$0.917_{\pm 0.074}$
	EquIN10	<b><math>0.003_{\pm 0.001}</math></b>	<b><math>0.910_{\pm 0.011}</math></b>
SOLIDS	ENR	$0.037_{\pm 0.038}$	$0.817_{\pm 0.085}$
	Baseline	-	$0.123_{\pm 0.007}$
	EquIN12	-	$0.126_{\pm 0.004}$
	EquIN24	-	$0.139_{\pm 0.056}$
	EquIN60	-	$0.596_{\pm 0.106}$
	EquIN80	-	<b><math>0.795_{\pm 0.230}</math></b>
ENR	-	$0.772_{\pm 0.095}$	

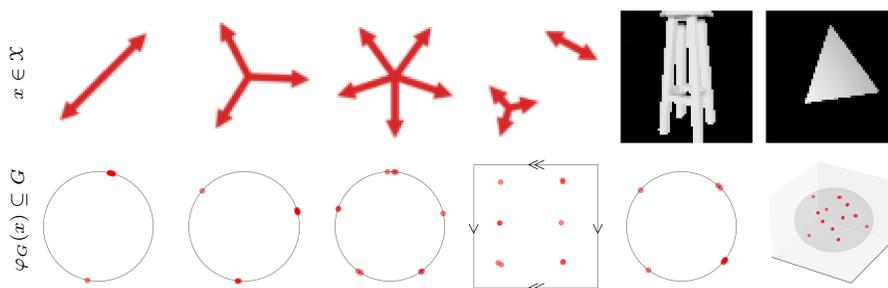


Fig. 4: Visualization of datapoints  $x$  and the corresponding predicted (coset of the) stabilizer  $\varphi_G(x)$ . For DOUBLE ARROWS, the torus  $G = \text{SO}(2) \times \text{SO}(2)$  is visualized as an identified square. For the tetrahedron from SOLIDS,  $G$  is visualized as a projective space  $\mathbb{RP}^3 \simeq \text{SO}(3)$ .

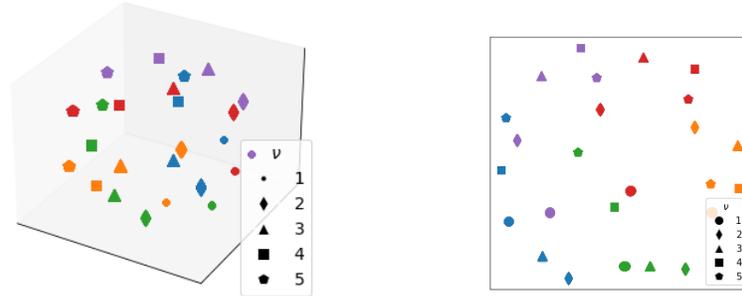


Fig. 5: Embeddings  $\varphi_{\mathcal{O}}(x) \in \mathcal{Z}_{\mathcal{O}} \subseteq \mathbb{R}^3$  for  $x$  in COLORED ARROWS. Each symbol represents the ground-truth cardinality  $\nu = |G_x|$  of the stabilizer while the color of the symbol represents the corresponding color of the arrow (left). The same embeddings are projected onto  $\mathbb{R}^2$  via principal component analysis (right).

number of stabilizers as expected. Figure 5 displays the output of  $\varphi_{\mathcal{O}}$  for data from COLORED ARROWS. The orbits are correctly separated in  $\mathcal{Z}_{\mathcal{O}}$ . Therefore, the model is able to distinguish data due to variability in the number  $\nu$  of arrows as well as in their color.

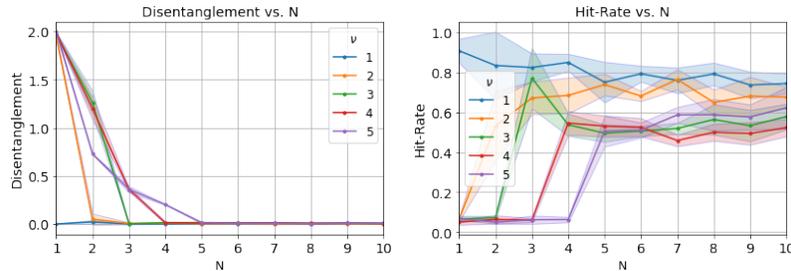


Fig. 6: Disentanglement and hit-rate for models trained with different values of  $N$ . Each line in the plot represents the results of a model trained on a dataset with a single orbit whose stabilizer has cardinality  $\nu$ . The plots show the mean and standard deviation across five repetitions.

## 5.5 Hyperparameter Analysis

For our last experiment, we investigate the effects of the hyperparameters  $N$  and  $\lambda$  when training EquIN on datasets with different numbers of stabilizers.

First, we show that a value of  $N$  larger than the cardinality of the stabilizers is necessary to achieve good values of disentanglement, and hit-rate for datasets with non-free group action, see Figure 6. However, large values of  $N$  can result in

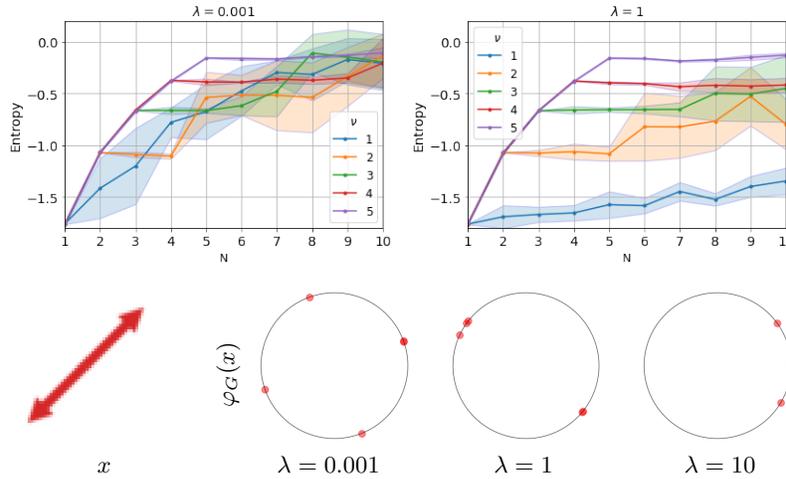


Fig. 7: Discrete entropy for models trained on the arrows dataset with different cardinalities of stabilizer  $\nu$  and two distinct values of  $\lambda$  (top row). Example embeddings  $\varphi_G(x)$  obtained for a datapoint  $x$  with two stabilizers obtained with models using  $\lambda \in \{0.001, 1, 10\}$  (bottom row).

non-collapsing embeddings  $\varphi_G$  corresponding to non-minimal cosets of the stabilizers. In these cases, the regularization term of Equation 9 and its corresponding weight  $\lambda$  plays an important role.

The bottom row of Figure 7 shows the embeddings  $\varphi_G(x)$  learnt for a datapoint  $x \in \mathcal{X}$  with stabilizer  $G_x \simeq C_2$  of cardinality two. The plots show how for low values of  $\lambda$ , the network converges to a non-minimal set. When an optimal value is chosen, such as  $\lambda = 1$ , the embeddings obtained with  $\varphi_G$  collapse to a set with the same cardinality as the stabilizers. If  $\lambda$  is too large, the embeddings tend to degenerate and collapse to a single point.

If the value of  $\lambda$  is too small, the discrete entropy of the learnt embeddings is not restricted. It continues to increase even if the number of embeddings matches the correct number of stabilizers. When an appropriate value of  $\lambda$  is chosen, the entropy becomes more stable as the embeddings have converged to the correct cardinality.

The plots in Figure 8 show the inverse relationship between  $\lambda$  and the entropy of the encoder  $\varphi_G$  that describes the collapse of the embeddings. The collapse of the embeddings also results in a lower performance of disentanglement and hit-rate by the models as seen for higher values of  $\lambda > 1$ . Throughout the experiments, we fix the value of  $\lambda = 1$  except for SOLIDS where a value of  $\lambda = 10$  was chosen since the number  $N$  used is larger.

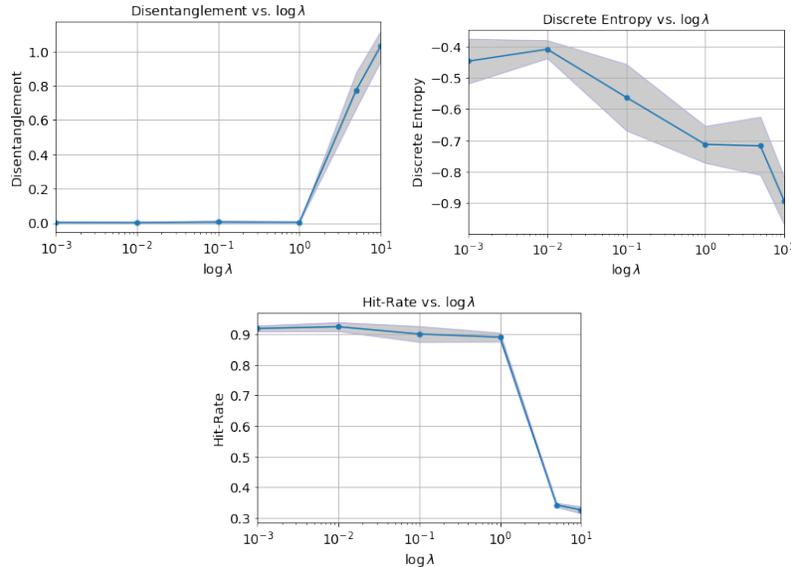


Fig. 8: Disentanglement, discrete entropy and hit-rate for models trained with different values of  $\lambda$  and fixed  $N = 5$ . The training dataset corresponds to the rotating arrows with  $\nu \in \{1, 2, 3, 4, 5\}$ . Each line shows the mean and standard deviation across five repetitions.

## 6 Conclusions and Future Work

In this work, we introduced EquIN, a method for learning equivariant representations for possibly non-free group actions. We discussed the theoretical foundations and empirically investigated the method on images with rotational symmetries. We showed that our model can capture the cosets of the group stabilizers and separate the information characterizing multiple orbits.

EquIN relies on the assumption that the stabilizers of the group action are finite. However, non-discrete stabilizer subgroups sometimes occur in practice, e.g., in continuous symmetrical objects such as cones, cylinders or spheres. Therefore, an interesting future direction is designing an equivariant representation learner suitable for group actions with non-discrete stabilizers.

## Acknowledgements

This work was supported by the Swedish Research Council, the Knut and Alice Wallenberg Foundation and the European Research Council (ERC-BIRD-884807). This work has also received funding from the NWO-TTW Programme “Efficient Deep Learning” (EDL) P16-25.

## Ethical Statement

The work presented in this paper consists of a theoretical and practical analysis on learning representations that capture the information about symmetry transformations observed in data. Due to the nature of this work as fundamental research, it is challenging to determine any direct adverse ethical implications that might arise. However, we think that any possible ethical implications of these ideas would be a consequence of the possible applications to augmented reality, object recognition, or reinforcement learning among others. The datasets used in this work consist of procedurally generated images with no personal or sensitive information.

## References

1. Ahuja, K., Hartford, J., Bengio, Y.: Properties from mechanisms: An equivariance perspective on identifiable representation learning. In: International Conference on Learning Representations (2022)
2. Bengio, Y., Courville, A., Vincent, P.: Representation learning: A review and new perspectives. In: IEEE Transactions on Pattern Analysis and Machine Intelligence (2013)
3. Caselles-Dupré, H., Garcia-Ortiz, M., Filliat, D.: Symmetry-based disentangled representation learning requires interaction with environments. In: Advances in Neural Information Processing Systems (2019)
4. Chen, R.T., Li, X., Grosse, R.B., Duvenaud, D.K.: Isolating sources of disentanglement in variational autoencoders. In: Advances in Neural Information Processing Systems (2018)
5. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International Conference on Machine Learning (2020)
6. Cohen, T., Welling, M.: Learning the irreducible representations of commutative Lie groups. In: International Conference on Machine Learning (2014)
7. Dupont, E., Martin, M.B., Colburn, A., Sankar, A., Susskind, J., Shan, Q.: Equivariant neural rendering. In: International Conference on Machine Learning (2020)
8. Falorsi, L., de Haan, P., Davidson, T.R., Cao, N.D., Weiler, M., Forré, P., Cohen, T.S.: Explorations in homeomorphic variational auto-encoding. In: ICML18 Workshop on Theoretical Foundations and Applications of Deep Generative Models (2018)
9. Guo, X., Zhu, E., Liu, X., Yin, J.: Affine equivariant autoencoder. In: International Joint Conference on Artificial Intelligence (2019)
10. Ha, D., Schmidhuber, J.: World models. arXiv preprint (2018)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
12. Higgins, I., Amos, D., Pfau, D., Racaniere, S., Matthey, L., Rezende, D., Lerchner, A.: Towards a definition of disentangled representations. arXiv preprint (2018)
13. Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., Lerchner, A.: Beta-vae: Learning basic visual concepts with a constrained variational framework. In: International Conference on Learning Representations (2017)

14. Higgins, I., Racanière, S., Rezende, D.: Symmetry-based representations for artificial and biological general intelligence. *Frontiers in Computational Neuroscience* (2022)
15. Hinton, G.E., Krizhevsky, A., Wang, S.D.: Transforming auto-encoders. In: *Artificial Neural Networks and Machine Learning*. Springer
16. Kipf, T., van der Pol, E., Welling, M.: Contrastive learning of structured world models. In: *International Conference on Learning Representations* (2020)
17. Liang, K.K.: Efficient conversion from rotating matrix to rotation axis and angle by extending rodrigues' formula. *arXiv preprint* (2018)
18. Locatello, F., Bauer, S., Lucic, M., Raetsch, G., Gelly, S., Schölkopf, B., Bachem, O.: Challenging common assumptions in the unsupervised learning of disentangled representations. In: *International Conference on Machine Learning* (2019)
19. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: *International Conference on Learning Representations* (2 2019)
20. Marchand, E., Uchiyama, H., Spindler, F.: Pose estimation for augmented reality: A hands-on survey. *IEEE Transactions on Visualization and Computer Graphics* (2016)
21. Marchetti, G.L., Tegnér, G., Varava, A., Kragic, D.: Equivariant Representation Learning via Class-Pose Decomposition. *arXiv preprint* (2022)
22. Murphy, K., Esteves, C., Jampani, V., Ramalingam, S., Makadia, A.: Implicit representation of probability distributions on the rotation manifold. In: *International Conference on Machine Learning* (2021)
23. Painter, M., Hare, J., Prügel-Bennett, A.: Linear disentangled representations and unsupervised action estimation. In: *Advances in Neural Information Processing Systems* (2020)
24. Park, J.Y., Biza, O., Zhao, L., van de Meent, J.W., Walters, R.: Learning symmetric embeddings for equivariant world models. *International Conference on Machine Learning* (2022)
25. van der Pol, E., Kipf, T., Oliehoek, F.A., Welling, M.: Plannable approximations to mdp homomorphisms: Equivariance under actions. In: *International Conference on Autonomous Agents and Multi-Agent Systems* (2020)
26. Quessard, R., Barrett, T.D., Clements, W.R.: Learning disentangled representations and group structure of dynamical environments. In: *Advances in Neural Information Processing Systems* (2020)
27. Rotman, J.J.: *An introduction to the theory of groups*, vol. 148. Springer Science & Business Media (2012)
28. Tonnaer, L., Perez Rey, L.A., Menkovski, V., Holenderski, M., Portegies, J.: Quantifying and Learning Linear Symmetry-Based Disentanglement. In: *International Conference on Machine Learning* (2022)
29. Worrall, D.E., Garbin, S.J., Turmukhambetov, D., Brostow, G.J.: Interpretable transformations with encoder-decoder networks. In: *International Conference on Computer Vision* (2017)
30. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2015)