

Learning Single/Multi-Attribute of Object with Symmetry and Group

Yong-Lu Li, Yue Xu, Xinyu Xu, Xiaohan Mao, Cewu Lu, *Member, IEEE*

Abstract—Attributes and objects can compose diverse compositions. To model the compositional nature of these concepts, it is a good choice to learn them as transformations, e.g., coupling and decoupling. However, complex transformations need to satisfy specific principles to guarantee rationality. Here, we first propose a previously ignored principle of attribute-object transformation: **Symmetry**. For example, coupling `peeled-apple` with attribute `peeled` should result in `peeled-apple`, and decoupling `peeled` from `apple` should still output `apple`. Incorporating the symmetry, we propose a transformation framework inspired by group theory, i.e., **SymNet**. It consists of two modules: Coupling Network and Decoupling Network. We adopt deep neural networks to implement **SymNet** and train it in an end-to-end paradigm with the group axioms and symmetry as objectives. Then, we propose a Relative Moving Distance (RMD) based method to utilize the attribute change instead of the attribute pattern itself to classify attributes. Besides the compositions of single-attribute and object, our RMD is also suitable for complex compositions of multiple attributes and objects when incorporating attribute correlations. **SymNet** can be utilized for attribute learning, compositional zero-shot learning and outperforms the state-of-the-art on four widely-used benchmarks. Code is at <https://github.com/DirtyHarryLYL/SymNet>.

Index Terms—Attribute-Object Composition, Compositional Zero-shot Learning, Single/Multi-Attribute, Symmetry, Group Axioms.



1 INTRODUCTION

ATTRIBUTES describe the properties of generic objects, e.g., material, color, weight, etc. Understanding the attributes would directly facilitate many tasks that require deep semantics, such as scene graph generation [9], object perception [8], human-object interaction detection [30], [31], [32], [48], [49]. As side information, attributes can also be employed in zero-shot learning [10], [11], [12], [14], [15], [29].

Going along with the road of conventional classification, some works [12], [13], [19], [24] address attribute recognition with discriminative models for objects but achieve poor performance. This is because attributes cannot be well expressed independently of the context [1], [2] (Fig. 1(a)). Subsequently, researchers rethink the nature of attributes and treat them as linear operations [2] to operate these general concepts: “adding” attribute to objects (coupling) or “removing” attribute from objects (decoupling). Though such insight has promoted this field, the “add-remove” system is not complete and lacks an axiomatics foundation to satisfy the specific principles of nature.

In this paper, we rethink the *physical* and *linguistic* properties of attribute-object, and propose a previously ignored but important principle of attribute-object transformations: **symmetry**, which would promote attribute-object learning. Symmetry depicts the invariance under transformations, e.g., a circle has rotational symmetry under the rotation without changing its appearance. The transformation that “adding” or “removing” attributes should also satisfy the symmetry: an object should remain unchanged if we “add”

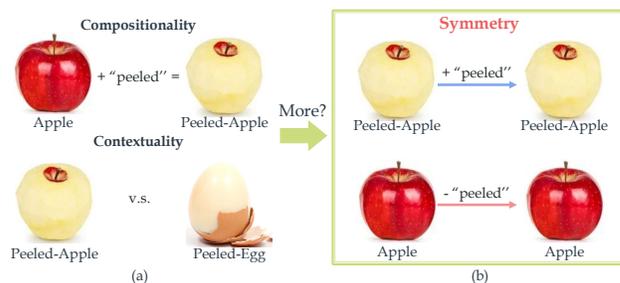


Fig. 1: Except for the compositionality and contextuality, attribute-object compositions also have the *symmetry* property. For instance, a `peeled-apple` should not change after “adding” the `peeled` attribute. Similarly, an `apple` should keep the same after “removing” the `peeled` attribute because it does not have it.

an attribute it already has, or “remove” an attribute it does not have. For instance, a `peeled-apple` keeps invariant if we “add” attribute `peeled` upon it. Similarly, “removing” attribute `peeled` from `apple` would still result in `apple`.

As shown in Fig. 1(b), except for the compositionality and contextuality, the symmetry should also be satisfied to guarantee rationality. Given this, we first introduce the symmetry and propose **SymNet** to depict it. In this work, we aim to bridge attribute-object learning and group theory. The elegant properties of groups would largely help in a more principled way, given its theoretical potential. Thus, to cover the principles existing in transformations theoretically, we borrow the principles from group theory to model symmetry. In detail, we define three transformations {“keep”, “add”, “remove”} and an operation to perform three transformations upon objects, in other words, to construct a “group”. To implement these, **SymNet** adopts Coupling Network (CoN) and Decoupling Network (DecoN) to perform coupling/adding and decoupling/removing. To meet the fundamental requirements of group theory, *symmetry* and the group axioms *closure*, *associativity*, *identity*

- Yong-Lu Li, Yue Xu, Xinyu Xu, Xiaohan Mao are with the Department of Electrical and Computer Engineering, Shanghai Jiao Tong University, Shanghai, 200240, China. E-mail: {yonglu_li, silixuyue, xuxinyu2000, mxh1999}@sjtu.edu.cn.
- Cewu Lu is the corresponding author, member of Qing Yuan Research Institute and MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, China. E-mail: lucewu@sjtu.edu.cn.

element, *invertibility element* are all implemented as the learning objectives of SymNet. Naturally, SymNet considers the compositionality and contextuality during the coupling and decoupling of various attributes and objects. All the above principles will be learned under a unified model in an end-to-end paradigm.

With symmetry learning, we can apply SymNet to address the Compositional Zero-Shot Learning (CZSL) task, whose target is to classify unseen compositions composed of seen attributes and objects. We adopt a novel attribute recognition paradigm, **Relative Moving Distance (RMD)** (Fig. 2). That said, given a specific attribute, an object would be manipulated by the “add” and “remove” transformations parallelly in *latent* space. Then, we can discriminate the existence of an attribute when transformations meet the symmetry principle (Sec. 3.4.1): if the input object has this attribute, the output after addition should be close to the original input object, and the object after removal should be far from the input. Contrarily, if the object does not have this attribute, the object after removal should be closer to the input than the object after addition. So attribute classification can be accomplished concurrently by comparing the relative *moving* distances between the input and two outputs.

In CZSL, the composition consists of a single attribute and an object [10], [11]. However, in practical application, an object usually has multiple attributes simultaneously [12], [15]. Thus, multi-attribute recognition has greater practical significance. However, under the multi-attribute setting, the attribute **correlation** would complicate the RMD principles. For example, when removing attribute *fresh* from objects, a {*green, juicy*} object is more likely to **change more** than a {*black, hard*} object in latent space. The reason is: though neither object has attribute *fresh*, attributes {*green, juicy*} are more closely related to *fresh* than {*black, hard*}. Further experiments also show that vanilla RMD for single-attribute scenarios fails to model multi-attribute correlation. Therefore, we further incorporate attribute correlation into RMD principles to adapt to the multi-attribute setting (Sec. 3.4.2).

With RMD, we can utilize the robust *attribute change* to classify attributes, instead of only relying on the dramatically unstable *visual attribute patterns*. Extensive experiments show that our method achieves significant improvements on both single- and multi-attribute learning benchmarks [10], [11], [12], [15]. The main contributions of this work are: 1) We propose a novel property of attribute-object composition transformation: symmetry, and design a framework inspired by group theory to learn it under the supervision of group axioms. 2) Based on symmetry learning, we propose a novel method to infer attributes based on relative moving distance (RMD). 3) We propose the corresponding RMD constraints to guide the learning for both single- and multi-attribute settings. 4) Substantial improvements are achieved in attribute recognition and CZSL tasks.

2 RELATED WORK

Visual Attribute. The visual attribute was introduced into computer vision to reduce the gap between visual patterns and object concepts, such as reducing the difficulty in object recognition [15] or acting as an intermediate representation

for zero-shot learning [13], [35]. After that, attribute has been widely applied in recognition of face [22], people [21], pedestrian [20] or action [23], zero-shot learning [12], [29] and so on. Therefore, attribute recognition is a fundamental problem to promote visual concept understanding.

The typical approach for attribute recognition is to train a multi-attribute discriminative model same as object classification [12], [13], [19], [24]. It ignores the intrinsic properties of attributes, such as compositionality and contextuality. Farhadi et al. [15] propose a visual feature selection method to recognize the attributes under the consideration of cross-category generalization. Gan et al. [37] further enhance the generalization by integrating kernel alignment with distributional variance. Liang et al. [39] think visual attributes are class-sensitive and utilize category information to predict attributes in a unified manner. Attribute correlation is essential information and gets explicitly considered in some works [36], [38], [40], [41]. Choi et al. [36] propose a hyper-graph framework to learn the semantic attributes correlation and apply it to scene recognition. Hand and Chellappa [41] designs a multi-task deep neural network with an auxiliary relation network for attribute prediction. Following these works, our SymNet explores attribute correlation from a novel perspective of attribute transformation. Some works [38], [40] apply automatic learning techniques to design deep neural networks for multi-task attribute learning automatically. Tang et al. [45] improve attribute recognition with weakly supervised attribute-specific localization. Later, some works start to consider the intrinsic properties by exploiting the attribute-object correlation [25], [27], [28]. Considering the contextuality of attributes, Nagarajan et al. [2] regard attributes as linear transformations operated upon object embeddings, and Misra et al. [1] map the attributes into model weight space to attain better representations. Adversarial learning is employed to model the discrepancy and correlations among attributes and objects. Yang et al. [44] propose a hierarchical feature embedding framework with inter-class and intra-class relations. Li et al. [46] propose a structural attribute learning framework to extract domain-invariant attribute features. Moreover, multi-attribute compositions can also be used to describe objects in few-shot recognition [3]. Different from our attribute correlation constraint (Sec. 3.4.2), Tokmakov et al. [3] adopts an orthogonal constraint to deal with attribute correlation. In Suppl Sec. 4.3, we compare two constraints and illustrate their respective advantages.

Compositional Zero-Shot Learning. CZSL is a crossing field of compositional learning and zero-shot learning. In the CZSL setting, test compositions are unseen during training, while each component is seen in both the train and test sets. Chen et al. [25] construct linear classifiers for unseen compositions with tensor completion of weight vectors. Misra et al. [1] consider that the model space is more smooth, thus project the attributes or objects into it by training binary linear SVMs for the corresponding components. For CZSL, it composes attribute and object embeddings in model space as composition representation. Wang et al. [5] address the attribute-object compositional problem via conditional embedding modification, which relies on attribute word embedding [6] transformation. Nan et al. [4] map the image features and word vectors [16] into embedding

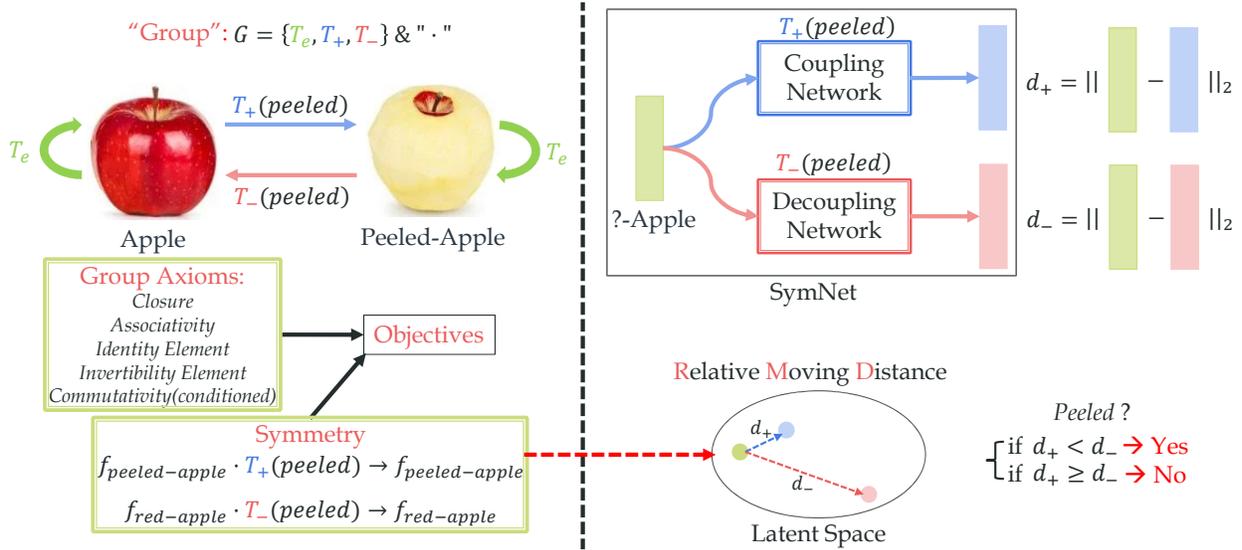


Fig. 2: Overview of our proposed method. We construct a “group” to learn the symmetry and operate the composition learning. The attribute transformations are implemented as coupling and decoupling networks and constrained by symmetry and group axiom objectives. Then relative moving distance based paradigm is applied in attribute classification.

space with the reconstruction constraint. Nagarajan et al. [2] regard attributes as linear operations for object embedding and map the image features and transformed object embeddings into a shared latent space. However, linear and explicit matrix transformation may be insufficient to represent various attribute concepts of different complexity, e.g., representing “red” and “broken” as matrices with the same capacity. Very recently, Naeem et al. [47] use a graph to learn the dependency and relevance between attributes, objects, and compositions. Previous methods usually ignored or incompletely considered the natural principles within the coupling and decoupling of attributes and objects. Hence, we propose a unified framework inspired by group theory to learn these essential principles such as symmetry.

3 APPROACH

Fig. 2 gives an overview of our approach. Our goal is to learn the symmetry within attribute-object compositions. Thus we can utilize it to obtain a deeper understanding of attribute-object, e.g., to address the CZSL task [10], [11]. To learn the symmetry in transformations, we need a comprehensive framework to cover all principles. Inspired by the group theory, we define a unified model named SymNet.

We define $G = \{T_e, T_+, T_-\}$ containing identity (“keep”), coupling (“add”), and decoupling (“remove”) transformations (Sec. 3.1) for each attribute and utilize Deep Neural Networks to implement them (Sec. 3.2). It is natural to adopt group theory as the close associations between symmetry and group to depict symmetry theoretically. As a group should satisfy the axioms, i.e., *closure*, *associativity*, *identity element*, *invertibility element*, we construct the learning objectives based on these axioms to train the transformations (Sec. 3.3). In addition, SymNet satisfies the *commutativity* under conditions. With these constraints, we can naturally guarantee compositionality and contextuality. *Symmetry* allows us to use a novel method, relative moving distance, to identify whether an object has a certain attribute with T_+ and T_- (Sec. 3.4) for CZSL (Sec. 3.5).

3.1 Group Definition

To depict the symmetry, we need first to define the transformations. Naturally, we need two reciprocal transformations to “add” and “remove” the attributes. Further, we need an axiomatic system to restrain the transformations and keep the rationality. Thus, we define three transformations $G = \{T_e, T_+, T_-\}$ and an operation “.”. In practice, it is difficult to strictly follow the theory considering the physical and linguistic truth. For example, the operation between attribute transformations “peeled · broken” is odd. Thus, our “operation” is defined to be operated upon object only.

Definition 1. Identity transformation T_e keep the attributes of object. Coupling transformation T_+ couples a specific attribute with an object. Decoupling transformation T_- decouples a specific attribute from an object.

Definition 2. Operation “.” performs transformations $\{T_e, T_+, T_-\}$ upon an object. Notably, operation “.” is not the dot product and we use this notation to maintain consistency with group theory.

More formally, for object $o \in \mathcal{O}$ and attribute $a^i, a^j \in \mathcal{A}$, $a^i \neq a^j$, where \mathcal{O} denotes object set and \mathcal{A} denotes attribute set, operation “.” performs transformations in G upon an object/image embedding:

$$\begin{aligned} f_o^i \cdot T_+(a^j) &= f_o^{ij}, \\ f_o^{ij} \cdot T_-(a^j) &= f_o^i, \\ f_o^i \cdot T_e &= f_o^i, \end{aligned} \quad (1)$$

where f_o^i means o has one attribute a^i and f_o^{ij} means o has two attributes a^i, a^j . Here we do not sign a specific object category and use o for simplicity.

Definition 3. G has the **symmetry** property if and only if $\forall a^i, a^j \in \mathcal{A}, a^i \neq a^j$:

$$f_o^i = f_o^i \cdot T_+(a^i), \quad f_o^i = f_o^i \cdot T_-(a^j). \quad (2)$$

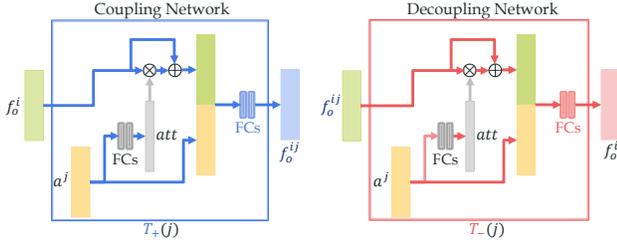


Fig. 3: The structure of CoN and DecoN. They take the attribute embedding to assign a specific attribute a^j . f_o^i, f_o^{ij} are the object embeddings extracted from ResNet [17]. The attribute embeddings are converted to attentions and applied on object embeddings, then compressed by MLPs to output transformed representations.

3.2 Group Implementation

In practice, when performing T_e upon f_o^i , we directly use f_o^i as the $f_o^i \cdot T_e$ to implement the identity transformation for simplicity. For T_+, T_- , we propose **SymNet** which consists of two modules: Coupling Network (**CoN**) and Decoupling Network (**DecoN**). CoN and DecoN have the same structure but independent weights and are trained with different tasks. As seen in Fig. 3, CoN and DecoN both take the object embedding f_o^i and attribute embedding a^j as inputs, and output the transformed object embedding. We use attribute category word vectors such as GloVe [16] or one-hot vectors to represent attributes. f_o^i is extracted by an ImageNet [18] pre-trained ResNet [17] from image I , i.e., $f_o^i = F_{res}(I)$.

Intuitively, attributes affect objects differently, e.g., “red” changes the color, “wooden” changes the texture. In CoN and DecoN, we use an *attribute-as-attention* strategy, i.e., using $att = g(a^j)$ as attention, where $g(\cdot)$ means two fully-connected (FC) and a Sigmoid layer. We concatenate $f_o^i \circ att + f_o^i$ with original a^j as the input and use two FC layers to perform the transformation.

3.3 Group Axioms as Objectives

According to group theory, SymNet should satisfy four group axioms: *closure*, *associativity*, *identity element*, and *invertibility*. Under certain conditions, attribute-object also satisfy *commutativity*. Besides, SymNet must obey the symmetry property of the attribute transformations.

In practice, we use Deep Neural Networks to implement transformations. Thus, we can construct training objectives to approach the theoretic transformations following the axioms. Considering the *actual characteristics* of attribute-object compositions, we slightly adjust the axioms to construct the objectives. Besides, there are two situations with different forms of axioms: 1) coupling or decoupling an attribute a^i that the object f_o^i already has, or 2) coupling or decoupling an attribute a^j that object f_o^i does not have.

Symmetry. First of all, SymNet should satisfy the symmetry property as depicted in Eq. 2, i.e., $f_o^i = f_o^i \cdot T_+(a^i), f_o^i = f_o^i \cdot T_-(a^j)$. The symmetry is essential to keep the semantic meaning during coupling and decoupling. For example, given a peeled-egg, adding the attribute peeled again should not change the object state. Similarly, a cup without attribute broken should remain unchanged after removing broken. Thus, we construct the **symmetry loss**:

$$\mathcal{L}_{sym} = \|f_o^i - f_o^i \cdot T_+(a^i)\|_2 + \|f_o^i - f_o^i \cdot T_-(a^j)\|_2, \quad (3)$$

where $a^i, a^j \in \mathcal{A}, i \neq j$. We use L_2 norm loss to measure the distance between two embeddings.

Closure. For all elements in set G , their operation results should also be in G . In SymNet, for the attribute a^i that f_o^i has, $f_o^i \cdot T_+(a^i) \cdot T_-(a^i)$ should be equal to $f_o^i \cdot T_-(a^i)$. For the attribute a^j that f_o^i does not have, $f_o^i \cdot T_-(a^j) \cdot T_+(a^j)$ should be equal to $f_o^i \cdot T_+(a^j)$. Thus, we construct:

$$\mathcal{L}_{clo} = \|f_o^i \cdot T_+(a^i) \cdot T_-(a^i) - f_o^i \cdot T_-(a^i)\|_2 + \|f_o^i \cdot T_-(a^j) \cdot T_+(a^j) - f_o^i \cdot T_+(a^j)\|_2. \quad (4)$$

Identity Element. The properties of identity element T_e are automatically satisfied since we implement T_e as a skip connection, i.e., $f_o^i \cdot T_*(a^i) \cdot T_e = f_o^i \cdot T_e \cdot T_*(a^i) = f_o^i \cdot T_*(a^i)$ where T_* denotes any element in G .

Invertibility Element. According to the definition, T_+ is the invertibility element of T_- , vice versa. For the attribute a^i that f_o^i has, $f_o^i \cdot T_-(a^i) \cdot T_+(a^i)$ should be equal to $f_o^i \cdot T_e = f_o^i$. For the attribute a^j that f_o^i does not have, $f_o^i \cdot T_+(a^j) \cdot T_-(a^j)$ should be equal to $f_o^i \cdot T_e = f_o^i$. Therefore, we have:

$$\mathcal{L}_{inv} = \|f_o^i \cdot T_+(a^j) \cdot T_-(a^j) - f_o^i \cdot T_e\|_2 + \|f_o^i \cdot T_-(a^i) \cdot T_+(a^i) - f_o^i \cdot T_e\|_2. \quad (5)$$

Associativity. Because of the practical physical meaning of attribute-object compositions, we only define the operation “.” that operates a transformation upon an object embedding in Sec. 3.1 and do not define the operation between transformations. Thus, we *relax* the constraint and do not construct an objective according to associativity in practice.

Commutativity. Because of the specialty of attribute-object, SymNet satisfies the *commutativity* when coupling and decoupling *multiple* attributes. Thus, $f_o^i \cdot T_+(a^i) \cdot T_-(a^j)$ should be equal to $f_o^i \cdot T_-(a^j) \cdot T_+(a^i)$:

$$\mathcal{L}_{com} = \|f_o^i \cdot T_+(a^i) \cdot T_-(a^j) - f_o^i \cdot T_-(a^j) \cdot T_+(a^i)\|_2. \quad (6)$$

Although the above definitions do not strictly follow the theory, the *loosely* conducted axiom objectives still contribute to the robustness and effectiveness a lot (ablation study in Sec. 4.9) and open the door to a more theoretical way.

Last but not least, CoN and DecoN need to keep the **semantic consistency**, i.e., before and after the transformation, the *object category* should not change. Hence, we use a cross-entropy loss \mathcal{L}_{cls}^o for the object recognition of the input and output embeddings of CoN and DecoN. In the same way, before and after coupling and decoupling, the *attribute changes* provide the attribute classification loss \mathcal{L}_{cls}^a . We use typical visual pattern-based classifiers consisting of FC layers for the object and attribute classification.

3.4 Relative Moving Distance

As shown in Fig. 4, we utilize the relative moving distance (RMD) based on the symmetry property to operate the attribute recognition. The implementations of RMD in single- and multi-attribute scenarios have minor differences.

3.4.1 Single-attribute RMD

Given an image embedding f_o^x of an object with an unknown attribute a^x , we input it to both CoN and DecoN with all kinds of attribute word embeddings $\{a^1, a^2, \dots, a^n\}$

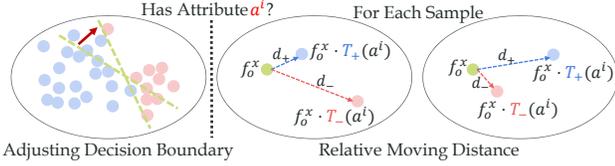


Fig. 4: Comparison between typical method and relative moving distance (RMD) based recognition. Previous methods mainly try to adjust the decision boundary. RMD based approach moves the embedding with T_+ and T_- and classifies by comparing their moving distances.

where n is the number of attributes. Two transformers would take attribute embeddings as conditions and operate coupling and decoupling *in parallel*, then output $2n$ transformed embeddings $\{f_o^x \cdot T_+(a^1), f_o^x \cdot T_+(a^2), \dots, f_o^x \cdot T_+(a^n)\}$ and $\{f_o^x \cdot T_-(a^1), f_o^x \cdot T_-(a^2), \dots, f_o^x \cdot T_-(a^n)\}$. We compute the distances between f_o^x and the transformed embeddings:

$$\begin{aligned} d_+^i &= \|f_o^x - f_o^x \cdot T_+(a^i)\|_2, \\ d_-^i &= \|f_o^x - f_o^x \cdot T_-(a^i)\|_2. \end{aligned} \quad (7)$$

To compare two distances, we define *relative moving distance* as $d^i = d_-^i - d_+^i$ and perform binary classification for each attribute (Fig. 4): 1) If $d^i \geq 0$, i.e., $f_o^x \cdot T_+(a^i)$ is closer to f_o^x than $f_o^x \cdot T_-(a^i)$, we tend to believe f_o^x has attribute a^i . 2) If $d^i < 0$, i.e., $f_o^x \cdot T_-(a^i)$ is closer, we tend to predict that f_o^x does not have attribute a^i . Previous zero/few-shot learning methods usually classify the instances via measuring the distance between the embedded instances and **fixed** points like prototype/label/centroid embeddings. Differently, **relative moving** distance compares the distance before and after applying the coupling and decoupling operations.

To enhance the RMD-based classification performance, we further use a triplet loss function. Let \mathcal{X} denote the attribute that f_o^x has, the loss can be described as:

$$\mathcal{L}_{tri} = \sum_i [d_+^i - d_-^i + \alpha]_+ + \sum_j [d_-^j - d_+^j + \alpha]_+, \quad (8)$$

where $\alpha=0.5$ is triplet margin, $[\cdot]_+ = \max(\cdot, 0)$. d_+^i should be less than d_-^i for the attribute that f_o^x has and greater than d_-^i for the attribute f_o^x does not have. The total loss is

$$\begin{aligned} \mathcal{L}_{total} &= \lambda_1 \mathcal{L}_{sym} + \lambda_2 \mathcal{L}_{axiom} + \lambda_3 \mathcal{L}_{cls}^a \\ &\quad + \lambda_4 \mathcal{L}_{cls}^o + \lambda_5 \mathcal{L}_{tri}, \end{aligned} \quad (9)$$

where $\mathcal{L}_{axiom} = \mathcal{L}_{clo} + \mathcal{L}_{inv} + \mathcal{L}_{com}$.

3.4.2 Multi-attribute RMD

When an object has multiple attributes, the RMD-based paradigm should be amended due to the attribute correlation. Fig. 5 illustrates the correlation matrices of attributes in aPY [15] and SUN [12], we can find that some highly correlated attributes exist. In transformations, the correlations between the removed attributes and existing attributes of an object should be considered. We modify the RMD triplet loss \mathcal{L}_{tri} from the following two aspects, as shown in Fig. 6. **Multi-attribute Symmetry Constraint.** Under the single-attribute setting, symmetry is strictly satisfied, i.e., the moving distance after adding an existing attribute (an object has) or removing a non-existing attribute (an object does not have) is *zero*. When it comes to the multi-attribute setting, this distance can be a small *non-zero* value due to the attribute correlation. For instance, if we **remove** attribute

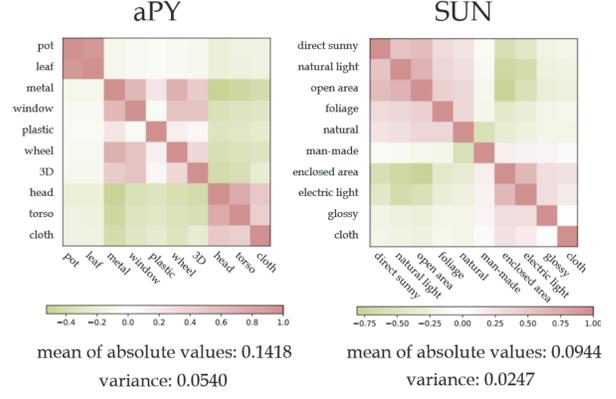


Fig. 5: Attribute correlation matrices of aPY [15] and SUN [12]. The correlation can be positive or negative ([-1,1], green to red). We calculate the mean and variance of the *absolute* correlation values, which show that aPY [15] contains *stronger* but more *variable* correlations.

metallic from a $\{\text{small, lustrous}\}$ object, it may lose part of attribute *lustrous*, leading to a small but non-zero moving distance. A similar example is shown in Fig. 6 (left). In the multi-attribute setting, for an object, we denote \mathcal{X} as its **all existing attributes**. Then, if attribute a^j (marked red in the figure) is more positively related to \mathcal{X} than a^i (marked blue), removing a^j would lead to a larger embedding moving distance. The above insight can be concluded into one simple sentence: for an object, removing an attribute that is more correlated to its existing attributes \mathcal{X} will make a more significant difference, and a less correlated attribute removing would make a minor difference.

Precisely, the correlations ([-1,1]) are measured via directly calculating the co-occurrence of attributes in the train set of multi-attribute benchmarks [12], [15]. For attributes a^i and a^j , let Y^i, Y^j denote their n -dimensional label vector among n total object instances in the train set, the **correlation coefficient** $\text{corr}(a^i, a^j)$ of a^i and a^j is measured as:

$$\text{corr}(a^i, a^j) = \frac{\text{cov}(Y^i, Y^j)}{\sqrt{\text{cov}(Y^i, Y^i)} \sqrt{\text{cov}(Y^j, Y^j)}}, \quad (10)$$

where $\text{cov}(Y^i, Y^j) = \frac{1}{n} (Y^i - \bar{Y}^i)^T (Y^j - \bar{Y}^j)$ and \bar{Y} denotes the mean value of vector Y . The definition of correlation can be generalized to that between attribute a^i and \mathcal{X} of an object without much effort:

$$\text{corr}(a^i, \mathcal{X}) = \sum_{a^j \in \mathcal{X}} \text{corr}(a^i, a^j). \quad (11)$$

We plot the moving distance after removing an attribute w.r.t. the correlation (between the removed attribute and \mathcal{X}) in Fig. 7. The moving distance and correlation of an ideal model should have a monotonically increasing relation (larger correlation causes a larger distance/change), while Fig. 7 shows that this relation is not well-learned by the SymNet with vanilla single-attribute symmetry since attributes are regarded as independent. Thus, we design an extra constraint in multi-attribute scenario. For an object with existing attributes \mathcal{X} , we randomly sample two **non-existing** attributes a^i and a^j (Fig. 6) and incorporate this property as a weighted triplet loss term:

$$\mathcal{L}_{tri}^{sym} = \left[(\text{corr}(a^i, \mathcal{X}) - \text{corr}(a^j, \mathcal{X})) (d_-^i - d_-^j) + \alpha \right]_+, \quad (12)$$

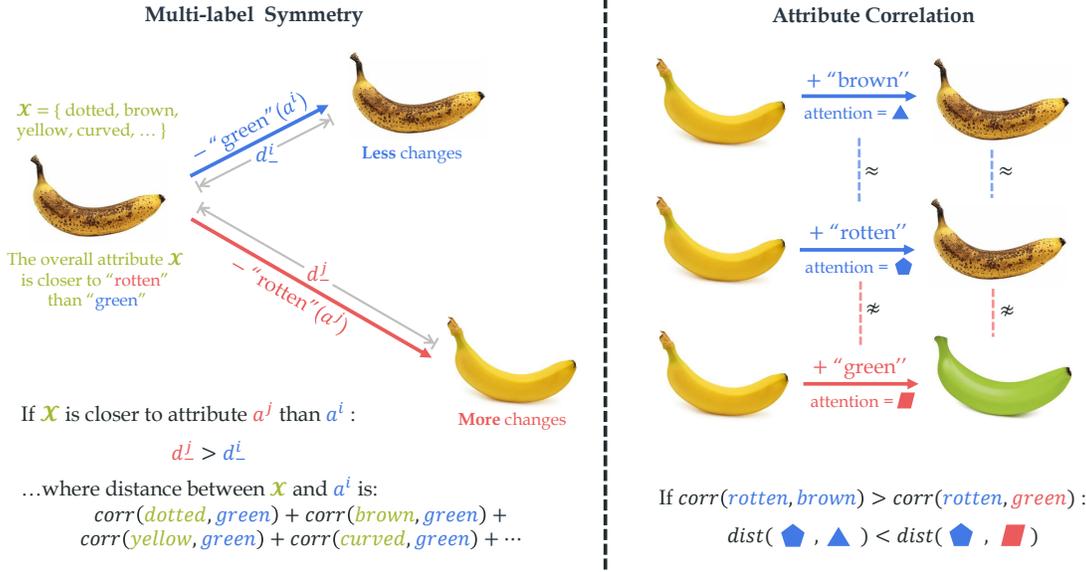


Fig. 6: Under multi-attribute/label setting, RMD need to consider more constraints since the existence of attribute correlation. (Left) For example, for a banana with existing attributes \mathcal{X} , since \mathcal{X} are closely related to rotten but have weak connection with green, the results differ when “removing” attribute rotten and green respectively. In single-attribute RMD, the results should be comparably close to the original banana, since this banana has neither attribute rotten nor green. But in multi-attribute RMD, the moving distance depends on the similarity of \mathcal{X} and the operated attribute. In general, removing an attribute which is more correlated to \mathcal{X} will make a larger difference, and the less correlated removing would make a minor difference. (Right) Meanwhile, for operated attributes, correlation would also affect the generated attentions: more correlated attributes should generate more similar attentions.

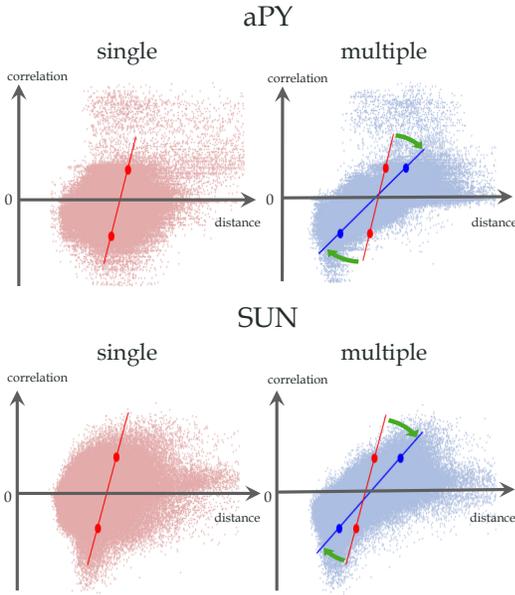


Fig. 7: The moving distances after attribute **removal** on multi-attribute benchmarks. The distance is always positive, and the correlation can be positive (> 0) or negative (< 0). According to our setting, **larger** positive correlation would generate **larger** moving distance in removal, e.g., removing attribute rotten in Fig. 6 (left). Oppositely, the **smaller** negative correlation could make the moving distance **smaller**, e.g., removing attribute green in Fig. 6 (left). The bold dots indicate the centroids of the dots with top-10% and down-10% correlation values. Thus, the line linked to two bold points should have a smaller slope ideally. The model is trained with single (left, red) or multiple (right, blue) RMD settings. We can find that the multi-attribute RMD generate more reasonable distances, i.e., has a smaller slope (blue lines).

where d^i_- is the moving distance as Eq. 7, α is the triplet margin, and $[\cdot]_+ = \max(\cdot, 0)$. Thus:

1) If $\text{corr}(a^i, \mathcal{X}) > \text{corr}(a^j, \mathcal{X})$, i.e., a^i is more correlated (larger correlation) with \mathcal{X} than a^j , then the triplet loss will

reduce d^j_- to make $d^j_- < d^i_-$.

2) If $\text{corr}(a^i, \mathcal{X}) < \text{corr}(a^j, \mathcal{X})$, i.e., a^j is more correlated with \mathcal{X} , the triplet loss will reduce d^i_- to achieve $d^j_- > d^i_-$.

Attribute Correlation Constraint. The correlation among attributes can also help attribute recognition since similar attributes similarly transform the object. For example, a $\{\text{muddy}, \text{dusty}\}$ object is very likely to be *dirty* since their high correlations. If a sub-optimal classifier gives high confidence to *muddy* and *dusty* but low confidence to *dirty*, incorporating the attribute correlation will help the classifier give consistent predictions on three similar attributes and therefore boost the performance.

The attention vectors from *attribute-as-attention* strategy imply how attributes manipulate object embeddings, so we implement the above correlation constraint as an auxiliary loss on generated attention. For each training object, we randomly sample three attributes a^i, a^j, a^k and compute their correlations $\text{corr}(a^i, a^j), \text{corr}(a^i, a^k)$ via Eq. 10. Then we construct an auxiliary triplet loss:

$$\begin{aligned} \mathcal{L}_{tri}^{corr} &= \left[\left(\text{corr}(a^i, a^j) - \text{corr}(a^i, a^k) \right) (d^i_+ - d^k_+) + \alpha \right]_+ \\ &+ \left[\left(\text{corr}(a^i, a^j) - \text{corr}(a^i, a^k) \right) (d^i_- - d^k_-) + \alpha \right]_+, \end{aligned} \quad (13)$$

where α is triplet margin, $[\cdot]_+ = \max(\cdot, 0)$. $d^i_+ = \text{dist}(\text{att}_+^i, \text{att}_+^j)$ is the distance between attribute attention vectors in $T_+(a^i)$ and $T_+(a^j)$, and $d^i_- = \text{dist}(\text{att}_-^i, \text{att}_-^j)$ is that in $T_-(a^i)$ and $T_-(a^j)$. The triplet loss is weighted by correlation difference $\text{corr}(a^i, a^j) - \text{corr}(a^i, a^k)$:

1) If $\text{corr}(a^i, a^j) > \text{corr}(a^i, a^k)$, i.e., a^i and a^j are more correlated, the loss will reduce d^i_+ and d^i_- to make the attentions of a^i and a^j more similar, so $d^i_+ < d^k_+, d^i_- < d^k_-$.

2) If $\text{corr}(a^i, a^j) < \text{corr}(a^i, a^k)$, i.e., a^i and a^k are more

correlated, the loss will reduce d_+^{ik} and d_-^{ik} to make the attentions of a^i and a^k more similar, so $d_+^{ij} > d_+^{ik}$, $d_-^{ij} > d_-^{ik}$.

The complete triplet loss in multi-attribute setting is

$$\mathcal{L}_{tri} = \sum_i^{\mathcal{X}} [d_+^i - d_-^i + \alpha]_+ + \sum_j^{A-\mathcal{X}} [d_-^j - d_+^j + \alpha]_+ \quad (14)$$

$$+ \lambda_6 \mathcal{L}_{tri}^{sym} + \lambda_7 \mathcal{L}_{tri}^{corr}.$$

Moreover, the total loss format is the same as Eq. 9.

3.4.3 Inference

In practice, for n attribute categories, we use RMDs $d = \{d^i\}_{i=1}^n$ as the attribute scores, i.e., $\mathcal{S}_a = \{\mathcal{S}_a^i\}_{i=1}^n = \{d^i\}_{i=1}^n$ and obtain attribute probability with Sigmoid function: $p_a^i = \text{Sigmoid}(\mathcal{S}_a^i)$. Notably, we also consider the scale and use a factor γ to adjust the score before Sigmoid. Our method can be operated in parallel, i.e., simultaneously computing the RMD values of n attributes. We input $[B, n, m]$ sized tensor where B is the mini-batch size and m is the object embedding size. CoN and DecoN would output two $[B, n, m]$ sized embeddings after transformation. Then we can compute RMDs $\{d^i\}_{i=1}^n$ simultaneously. Our method has approximately the same speed as a typical FC classifier. The inference speed from features to RMD is 41.0 FPS, and the FC classifier speed is 45.8 FPS. The gap can be further omitted if considering the overhead of the feature extractor.

3.5 Discussion: Composition Zero-Shot Learning

With robust and effective symmetry learning for attribute-object, we can further apply SymNet to CZSL [10], [11]. The goal of CZSL is to infer the unseen attribute-object pairs in the test set, i.e., a prediction is true positive if and only if both attribute and object classifications are accurate. The pair candidates are available during testing. Thus, the predictions of impossible pairs can be masked.

We propose a novel method to address this task based on RMD. With relative moving distance $d^i = d_-^i - d_+^i$, the attribute probability is computed as $p_a^i = \text{Sigmoid}(d^i)$. For the object category, we input the object embedding to 2-layer FC with Softmax to obtain the object scores $\mathcal{S}_o = \{\mathcal{S}_o^j\}_{j=1}^m$, where m is the number of object categories. The object probability $p_j = \text{Softmax}(\mathcal{S}_o^j)$ and $p_o = \{p_o^j\}_{j=1}^m$. We then use p_{ao}^{ij} to represent the probability of an attribute-object pair in the test set, which is composed of the i -th attribute category and j -th object. The pair probabilities are given by $p_{ao}^{ij} = p_a^i \times p_o^j$. The impossible compositions would be masked according to the benchmarks [10], [11].

4 EXPERIMENT

4.1 Data and Metrics

Our experiments are conducted on the following datasets (Suppl Sec. 2): MIT-States [10] and UT-Zappos50K [11] (single-attribute), aPY [15] and SUN [12] (multi-attribute).

In attribute recognition, for aPY and SUN, we report mAUC following previous methods such as [37], [38]. For MIT-States and UT-Zappos, we report Top-1 accuracy following [4]. The CZSL experiments are conducted on MIT-States and UT-Zappos, where the training and testing pairs are non-overlapping, i.e., the test set contains unseen

attribute-object pairs composed of seen attributes and objects. We report the Top-1, 2, 3 accuracies on the unseen test set. We also evaluate our model under the generalized CZSL setting of TMN [33], since the "open world" setting from [2] brings biases towards unseen pairs [34].

4.2 Baselines

We compare SymNet with previous state-of-the-arts (detailed in Suppl Sec. 2). For single-attribute learning and CZSL, we adopt the Visual Product, LabelEmbed (LE) [1], LabelEmbed Only Regression (LEOR) [1], LabelEmbed With Regression (LE+R) [1], LabelEmbed+ [2], AnalogousAttr [25], Red Wine [1], AttrOperator [2], TAFE-Net [5], GenModel [4], f-CLSWGAN [43], TMN [33], and Causal [42]. For multi-attribute learning, the ALE [35], HAP [36], UDICA/KDICA [37], UMF [39], AMF [41], FMT [40], GALM [38] are adopted.

4.3 Implementation Details

We use ImageNet pre-trained ResNet-18 [17] as backbone to extract features for MIT-States and UT-Zappos, ResNet-50 [17] for aPY and SUN. Especially, since images in aPY may have several instances, we use the feature after RoI-pooling [8]. We do not fine-tune it following previous methods. The 300-dimensional pre-trained GloVe [16] vectors are used as the word embeddings. Moreover, SymNet is trained with an SGD optimizer on a single NVIDIA GPU.

On MIT-States and UT-Zappos, the 512-dimensional ResNet-18 feature is first transformed to 300-dimensional by a single FC. On aPY and SUN, the 2048-dimensional ResNet-50 feature is compressed to 128-dimensional via an FC too. The main modules of our SymNet, CoN, and DecoN, have the same structures but independent weights as depicted in Fig. 3: two FC layers (512/300 on MIT-States UT-Zappos, 512/128 on aPY SUN) with Sigmoid convert the attribute embedding to attention with the same dimension for the object embedding. The attention is multiplied to the input object embedding and then gets summed with a shortcut of original object embedding. Next, the object embedding after the attention operation is concatenated to the attribute embedding and then compressed to the original dimension by the other two FC layers (256/128 on aPY, 1536/128 on SUN, 768/300 on MIT-States UT-Zappos). Each hidden FC in CoN and DecoN is followed by BatchNorm and ReLU.

On single-attribute benchmarks, for each training image, we randomly sample another image with the same object but a different attribute as the negative sample to compute L_{total} . (Sec. 3.3). On multi-attribute benchmarks, we compute $L_{sym}, L_{axiom}, L_{tri}, L_{cls}(\mathcal{L}_{cls}^a, \mathcal{L}_{cls}^o)$ with all attributes and randomly sample three different attributes to compute L_{tri}^{corr} for each object. Besides, we define the correlation of an attribute to an object as the sum of correlations between this attribute and \mathcal{X} of this object (Eq. 11). According to the order of correlation, we regard top-10% and last-10% attributes as the strongly related ones, middle-10% attributes as the neutral ones. One strongly related attribute and one neutral attribute would form a pair, and all such pairs are used to compute L_{tri}^{sym} . In the multi-attribute setting, coupling and decoupling already involve all attributes to implement the commutativity loss.

We use cross-validation to determine the hyper-parameters (Suppl Tab. 1). The weights on datasets are different as their different domains/ranges/scales, leading to distinct embedding spaces and different parameters.

4.4 Single-Attribute Learning

We first compare the attribute learning alone on two single-attribute benchmarks in Tab. 1. We reproduce the results of AttrOperator [2] with its open-sourced code. For all methods involved, the individual attribute and object accuracy do not consider the relations between attributes and objects. The object recognition module of our method is a simple 3-layer MLP classifier with the visual image features from the ResNet-18 backbone. SymNet outperforms previous methods by a large margin, i.e., 3.8% on MIT-States and 8.3% on UT-Zappos, which strongly verifies that our RMD-based attribute learning is particularly effective.

4.5 Multi-Attribute Learning

Next, we evaluate SymNet on multi-attribute learning benchmarks aPY [15] and SUN [12] in Tab. 2. SymNet still outperforms the state-of-the-art by 1.4% on aPY, 1.9% on SUN. The reason is the advantages of SymNet as a transformation-style framework considering the attribute correlation explicitly. Previous methods ignore the compositional nature of attributes and objects, e.g., UMF [39] directly uses the image and object representations in latent space to predict attributes. Relatively, SymNet captures how attribute interacts with object and model attribute-object transformation with **symmetry** and **group** principles. Thus, SymNet in the single-attribute setting directly outperforms the state-of-the-art on SUN. Besides, as seen in Fig. 5, attribute labels on aPY [15] are strongly correlated, but the single-attribute setting does not work well. Thus, the multi-attribute setting **explicitly** mines attribute correlation and uses it to regularize the attribute-attribute relationship in attribute-object transformation. Other methods like FMT [40], GALM [38] apply multi-tasking techniques to force the model to learn correlation automatically. Moreover, such an implicit mining approach leads to complex training, which is challenging to learn the correlation well.

Besides, in the single-attribute setting without both correlation-based losses (SymNet w/o \mathcal{L}_{tri}^{sym} & \mathcal{L}_{tri}^{corr}), the score is 82.1% on aPY and 88.1% on SUN, with a drop of 4.0% and 0.3% respectively. The reason is that, as revealed in Fig. 5, attribute correlation in aPY is much stronger. Thus, its performance gap between single- and multi-attribute settings is more significant. However, the attribute correlation in SUN is much softer, so the impact of multi-attribute constraint is slight, and the ablations without correlation losses result in comparable performances in two settings. These results show the effectiveness of our correlation-based loss and multi-attribute RMD and prove their general capability.

4.6 Compositional Zero-Shot Learning

To evaluate the symmetry learning in the compositional zero-shot task, we conduct experiments on widely-used benchmarks: MIT-States [10] and UT-Zappos [11].

Method	MIT-States	UT-Zappos
Visual Product [1]	14.7	24.9
AttrOperator [2]	14.6	29.7
GenModel [4]	15.1	18.4
SymNet	18.9	38.0

TABLE 1: Attribute learning results (accuracy, %) on single-attribute benchmarks.

Method	aPY	SUN
ALE [35]	69.2	74.5
HAP [36]	58.2	76.7
UDICA [37]	82.3	85.8
KDICA [37]	84.7	/
UMF [39]	79.7	80.5
AMT [41]	84.5	82.5
FMT [40]	70.5	75.5
GALM [38]	84.2	86.5
SymNet (single)	82.1	88.1
SymNet	86.1	88.4

TABLE 2: Attribute learning results (mAU, %) on multi-attribute benchmarks. SymNet (single) is the model without \mathcal{L}_{tri}^{sym} and \mathcal{L}_{tri}^{corr} .

The results of CZSL are shown in Tab. 3, where the first five rows are baselines from [1], [2] (the scores with * are reproduced by [2], the others are from [1]). SymNet outperforms all baselines on two benchmarks. Although we use a simple product to compose the attribute and object scores, we still achieve 2.1% and 3.8% improvements over the state-of-the-art [4] on two benchmarks, respectively. Most previous approaches do not surpass the Visual Product baseline on UT-Zappos, while ours outperforms by 2.2%.

In addition, our object classification performance on the two datasets are 28.8% and 65.4% respectively, which is comparable to AttrOperator [2] (20.5%, 67.5%) and GenModel [4] (27.7%, 68.1%). Accordingly, the main contribution of the CZSL improvement of SymNet comes from attribute learning rather than object recognition.

To further evaluate our SymNet, we additionally conduct the comparison on the *generalized* CZSL setting from recent state-of-the-art TMN [33] on the larger MIT-States [10] in Tab. 4. SymNet also outperforms previous methods significantly on all metrics. We notice that novel metrics on UT-Zappos [11] are proposed by a very recent work Causal [42], so we also conduct a test following its metrics. The results are shown in Tab. 5. Compared with TMN [33], SymNet shows its superiority on Seen, Harmonic, and AUC metrics. Furthermore, even compared with the very recent Causal [42] method, our SymNet still improves the Seen, Closed, and AUC metrics. All these results further strongly prove the efficacy of our method.

4.7 Application in Few-Shot Learning

SymNet is a generic feature extractor, and its features are rich in attribute semantics which can strengthen the object representations. As a verification, we conduct a few-shot recognition experiment on CUB-200-2011 [29] following the protocol of COMP [3]. We enhance the attribute representation for the downstream few-shot classification by concatenating the average CoN and DecoN transformed features over all attributes to the original ResNet feature, respectively. Setting details and supplementary results on SUN397 [12] are listed in Suppl Sec. 4.4.

We embed SymNet on the pre-trained ResNet-10 backbone and evaluate on both settings (whether use data augmentation during classifier training) following COMP [3].

Method	MIT-States			UT-Zappos		
	Top-1	Top-2	Top-3	Top-1	Top-2	Top-3
Visual Product [1]	9.8/13.9*	16.1	20.6	49.9*	/	/
LabelEmbed (LE) [1]	11.2/13.4*	17.6	22.4	25.8*	/	/
- LEOR [1]	4.5	6.2	11.8	/	/	/
- LE + R [1]	9.3	16.3	20.8	/	/	/
- LabelEmbed+ [2]	14.8*	/	/	37.4*	/	/
AnalogousAttr [25]	1.4	/	/	18.3	/	/
Red Wine [1]	13.1	21.2	27.6	40.3	/	/
AttrOperator [2]	14.2	19.6	25.1	46.2	56.6	69.2
TAFE-Net [5]	16.4	26.4	33.0	33.2	/	/
GenModel [4]	17.8	/	/	48.3	/	/
SymNet (Ours)	19.9	28.2	33.8	52.1	67.8	76.0

TABLE 3: CZSL results (top-k accuracy, %) on MIT-States [10] and UT-Zappos [11]. The scores with * mark are reproduced by [2] and the rest are reported in the original papers.

Method	Val AUC			Test AUC			Seen	Unseen	HM
	1	2	3	1	2	3			
AttrOperator [2]	2.5	6.2	10.1	1.6	4.7	7.6	14.3	17.4	9.9
Red Wine [1]	2.9	7.3	11.8	2.4	5.7	9.3	20.7	17.9	11.6
LabelEmbed+ [2]	3.0	7.6	12.2	2.0	5.6	9.4	15.0	20.1	10.7
f-CLSWGAN [43]	3.1	6.9	10.5	2.3	5.7	8.8	24.8	13.4	11.2
TMN [33]	3.5	8.1	12.4	2.9	7.1	11.5	20.2	20.1	13.0
SymNet	5.4	11.6	16.6	4.5	10.1	15.0	26.2	26.3	16.8

TABLE 4: Results of generalized CZSL on MIT-States [10] following [33]. All methods use ResNet-18 [17] as backbone.

As shown in Table 6, the SymNet features bring stable performance gain to the COMP [3] baselines on both novel and all categories in most situations, indicating that feature from SymNet attribute transformation is an effective supplementation of object representation. In addition, it also proves the generalization ability of SymNet cross object categories.

4.8 Visualization

To verify the robustness and principles in transformations, we use t-SNE [7] to visualize the attention vectors or image embeddings before or after transformations in Fig. 8.

Precisely, we first visualize the group axioms related transformations: 1) **Closure** is verified by comparing $\{f_o^i \cdot T_+(a^i) \cdot T_-(a^j)\}$ v.s. $\{f_o^i \cdot T_-(a^i)\}$ and $\{f_o^i \cdot T_-(a^j) \cdot T_+(a^j)\}$ v.s. $\{f_o^i \cdot T_+(a^j)\}$. 2) **Invertibility** is verified by comparing $\{f_o^i \cdot T_+(a^j) \cdot T_-(a^j)\}$ v.s. $\{f_o^i \cdot T_e\}$ and $\{f_o^i \cdot T_-(a^i) \cdot T_+(a^i)\}$ v.s. $\{f_o^i \cdot T_e\}$. 3) **Commutativity** is verified by comparing $\{f_o^i \cdot T_+(a^i) \cdot T_-(a^j)\}$ v.s. $\{f_o^i \cdot T_-(a^j) \cdot T_+(a^i)\}$. The results are shown in Fig. 8 (a,b). We observe that SymNet can robustly operate the transformations and the axiom objectives are well satisfied during embedding transformations.

Then, to verify the **symmetry** property, we visualize the sample embeddings in relative moving space in Fig. 8(c,d): 1) For the sample f_o^i which do not have attribute a^j , $f_o^i \cdot T_+(a^j)$ should be far from f_o^i . On the contrary, $f_o^i \cdot T_-(a^j)$ are relatively close to f_o^i because of the symmetry principle. 2) For the sample f_o^i with attribute a^i , $f_o^i \cdot T_+(a^i)$ should be close to f_o^i and $f_o^i \cdot T_-(a^i)$ should be far from f_o^i . We can also find that the relative moving distance rules are satisfied, i.e., the symmetry is well learned by our SymNet.

We also verify the properties in multi-attribute scenario in Fig. 8(e,f,g,h): 1) **Multi-Attribute Symmetry Constraint**: after removing an attribute which has *larger* positive correlation to an object, the moving distance would also be *larger*. In Fig. 8(e,f), the red dots are the positions of the original object embeddings. The others are the embeddings after the attribute removals. The dot color is related to the correlation

Model	Unseen	Seen	Harmonic	Closed	AUC
LabelEmbed+ [2]	16.2	53.0	24.7	59.3	22.9
AttrOperator [2]	25.5	37.9	27.9	54.0	22.1
TMN [33]	10.3	54.3	17.4	62.0	25.4
Causal [42]	28.0	37.0	30.6	58.6	26.4
SymNet	10.3	56.3	24.1	58.7	26.8

TABLE 5: Results of generalized CZSL on UT-Zappos [11] following [42]. All methods use ResNet-18 [17] as backbone. The results of LabelEmbed+, AttrOperator, and TMN are from [42].

Method	Novel			All		
	1-shot	2-shot	5-shot	1-shot	2-shot	5-shot
COMP [3]	52.5	63.6	73.8	62.6	68.4	74.0
COMP - SymNet	54.0	63.8	74.3	63.1	68.9	74.1
COMP w/ data aug [3]	53.6	64.8	74.6	63.1	69.2	74.5
COMP w/ data aug - SymNet	57.3	66.7	76.0	64.9	68.8	74.9

TABLE 6: Results of few-shot recognition on CUB-200-2011 [29].

value. Darker color denotes the removal of a *more correlated* attribute. We find that *larger* correlations would cause *larger* moving distances. This phenomenon depicts that the transformations are more principled and consistent with the practice with our multi-attribute symmetry constraint. 2) **Attention Correlation Constraint**: the attention vectors of *more correlated* attributes should be *closer*. In Fig. 8(g,h), this constraint indeed leads to better clusters. We can observe that our multi-attribute model can capture the attribute correlations and well-learn these two constraints. Besides, we report the **image retrieval** via SymNet in Suppl Sec. 3.

4.9 Ablation Study

To evaluate the components of our model, we compare the results of different model designs in Tab. 7 following the settings of [2]. We also conduct experiments of the generalized CZSL setting [33], COMP [3], and AttrOperator [2], which are attached in supplementary material.

(1) **Objectives**: to evaluate the objectives constructed from group axioms, the core principle symmetry, and attribute correlation in the multi-attribute scenario, we conduct experiments of removing objectives. In Tab. 7 and Suppl Tab. 2, SymNet shows obvious degradations without the constraints of these principles ($\mathcal{L}_{sym}, \mathcal{L}_{axiom}, \mathcal{L}_{cls}, \mathcal{L}_{tri}$). The degradations are in line with our assumption that a transformation framework that covers the essential principles can largely promote attribute learning.

1-a) Specifically, removing \mathcal{L}_{cls} leads to a more significant performance drop than other losses. Because \mathcal{L}_{cls} applied on the input/output **object embeddings** of CoN and DecoN can establish the basic distribution of object embeddings, thus keep the basic semantics of the object and attribute. While the other transformation-related losses (e.g., $\mathcal{L}_{axiom}, \mathcal{L}_{sym}, \mathcal{L}_{tri}$) ensure the attribute transformation rational and the object embedding moving to follow the theoretical guidance, thus can afford the robust classification via RMD. Removing \mathcal{L}_{cls} would destroy the basic semantic information within object embeddings (e.g., the relationship between attributes and objects, or the inter-attribute correlation), therefore, results in a larger performance drop. Similar phenomenon also presents in AttrOperator [2], that removing \mathcal{L}_{aux} (counterpart of our \mathcal{L}_{cls} , ensuring that the identity of the attribute/object is not lost during composition), leads to about relative 55% drop of h-mean.

1-b) We conduct an ablation of **only keeping \mathcal{L}_{cls} and \mathcal{L}_{tri}** . The top-1 accuracy drops 2.8%, 3.6%, 14.5%, and 3.5%

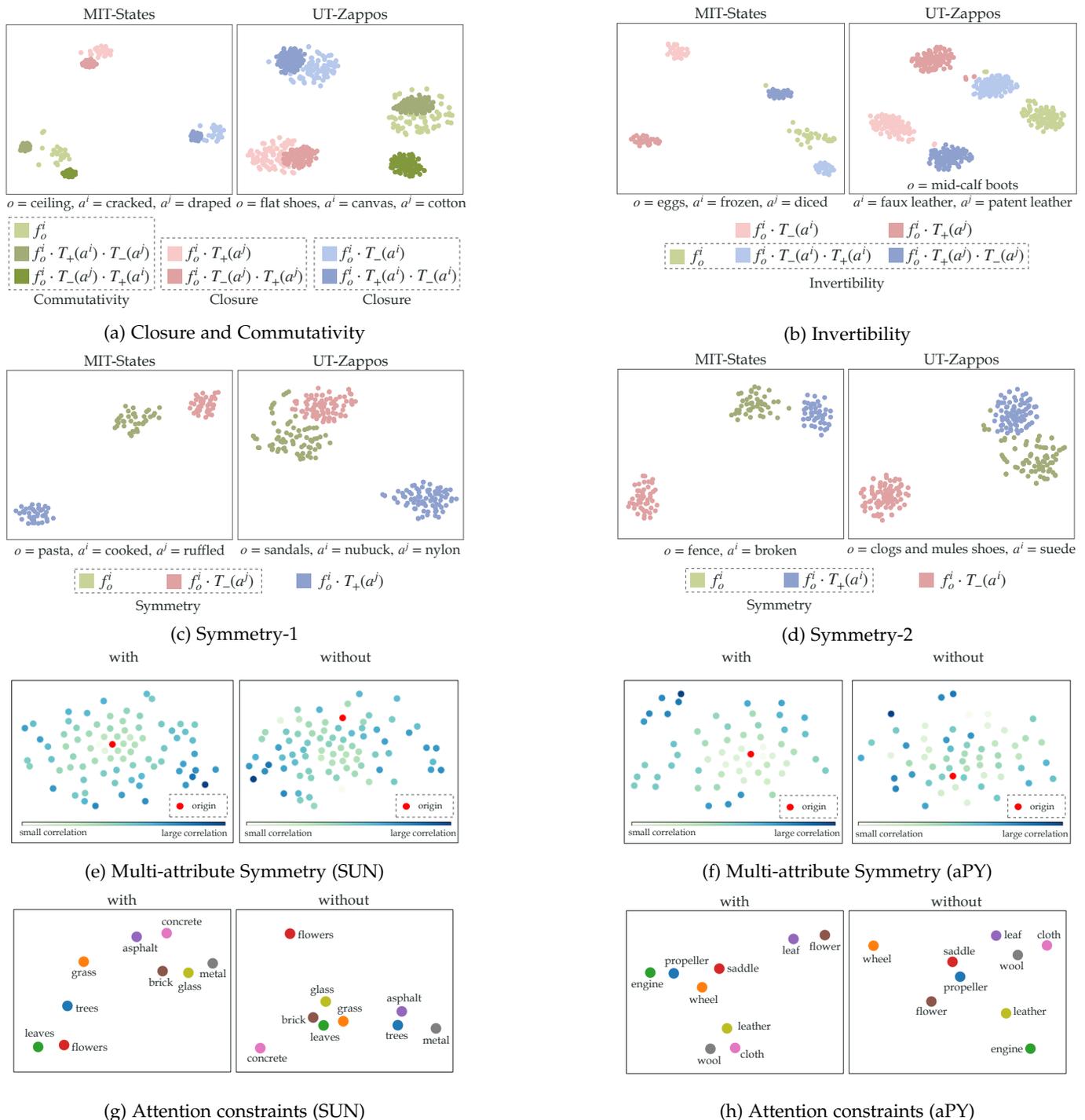


Fig. 8: Visualization of symmetry, group axioms and attention by t-SNE [7]. In (a)-(d), points with colors in the same dotted box should be close according to the corresponding learning principles. Especially, (e) and (f) shows the RMD property in the multi-attribute setting where the red dot is the original object embedding and the other dots are the embeddings after attribute removals. For example, if attribute a^i has a larger correlation to the existing attributes of this object, as $\text{corr}(a^i, \mathcal{X})$ is large, then the corresponding dot is in a darker color and has a larger distance after removal. On the contrary, the dot of removing attribute a^j with smaller $\text{corr}(a^j, \mathcal{X})$ is represented in a lighter color and closer to the origin. However, without the L_{trip}^{sym} , the distribution of attribute removal is much noisier. (g) and (f) are the distributions of attention representation. Attentions of more correlated attributes are closer in latent space, e.g., leaves and flowers, but this property cannot keep without L_{tri}^{corr} .

on four datasets in CZSL and the top-1 AUC drops 1.7 on MIT-States in generalized CZSL. This performance gap can then be filled by our proposed \mathcal{L}_{axiom} and \mathcal{L}_{sym} .

1-c) Though the overall score drop of “w/o \mathcal{L}_{cls} ” is larger than the transformation-related losses, the situations of specific attribute classes differ. In detail, we compare our best model to the model trained with only \mathcal{L}_{cls} & \mathcal{L}_{tri} . From

the results, we find that the symmetry and axiom constraints can facilitate the learning of **few-shot attributes**. Our full model outperforms the one with only \mathcal{L}_{cls} & \mathcal{L}_{tri} on fewer-shot attributes, such as *tight*, *bent*, *viscous*. The accuracy of samples of the thirty least frequent attributes increases by 4.3% with \mathcal{L}_{sym} & \mathcal{L}_{axiom} . While on the rest attributes with more samples, the full model only has 0.9% ac-

Method	MIT-States			UT-Zappos			aPY	SUN
	Top-1	Top-2	Top-3	Top-1	Top-2	Top-3		
SymNet	19.9	28.2	33.8	52.1	67.8	76.0	86.1	88.4
SymNet w/o \mathcal{L}_{sym}	18.3	27.5	33.4	51.1	67.0	76.0	65.7	84.0
SymNet w/o \mathcal{L}_{axiom}	16.9	25.5	30.9	47.6	65.4	73.6	83.6	87.9
SymNet w/o \mathcal{L}_{triv}	17.9	26.7	32.5	50.8	67.4	76.1	84.9	88.1
SymNet w/o \mathcal{L}_{com}	17.8	27.0	32.7	51.2	67.6	75.8	/	/
SymNet w/o \mathcal{L}_{clo}	18.0	27.0	32.8	51.1	67.2	76.0	84.1	88.1
SymNet w/o \mathcal{L}_{cls}	10.3	18.9	25.9	28.7	51.2	65.2	81.3	78.0
SymNet w/o \mathcal{L}_{tri}	17.8	26.8	32.6	49.2	65.3	74.2	83.5	88.0
SymNet w/o \mathcal{L}_{tri}^{sym}	/	/	/	/	/	/	84.0	88.1
SymNet w/o \mathcal{L}_{tri}^{corr}	/	/	/	/	/	/	84.5	88.1
SymNet w/o $\mathcal{L}_{tri}^{sym} \& \mathcal{L}_{tri}^{corr}$	/	/	/	/	/	/	82.1	88.1
SymNet w/o $\mathcal{L}_{sym} \& \mathcal{L}_{tri}$	17.7	27.0	33.0	50.1	66.1	75.6	63.1	85.3
SymNet w/o $\mathcal{L}_{tri} \& \mathcal{L}_{cls}$	10.5	19.4	26.7	28.6	51.4	65.6	80.1	70.9
SymNet w/o $\mathcal{L}_{sym} \& \mathcal{L}_{cls}$	9.3	17.0	22.7	27.4	48.2	64.1	66.4	78.0
SymNet only \mathcal{L}_{sym}	9.4	16.9	22.5	20.4	38.9	53.5	73.4	71.0
SymNet only $\mathcal{L}_{cls}, \mathcal{L}_{tri}$	17.1	26.1	31.7	48.5	65.7	73.9	71.6	84.9
SymNet w/o attention	18.0	26.9	32.7	48.5	65.0	75.6	84.8	88.2
SymNet tanh attention	16.9	25.0	30.8	42.0	59.0	69.0	84.1	88.0
SymNet L_1 dist.	7.1	11.2	14.3	37.5	53.3	62.3	82.0	87.6
SymNet Cos dist.	11.3	20.7	28.5	18.7	41.1	60.0	82.6	86.6

TABLE 7: Ablation studies on CZSL and multi-attribute learning.

curacy gain compared with the model with only $\mathcal{L}_{cls} \& \mathcal{L}_{tri}$. The reason may be that the RMD relying on symmetry and dynamic embedding moving ($\mathcal{L}_{sym} \& \mathcal{L}_{axiom}$) is less reliant on the training sample scale than canonical classification.

1-d) Noticeably, the difference of *datasets and metrics* also exert an influence. That said, the effects of losses vary on different datasets depending on the data scale, quality, and distribution. For example, comparing the single- and multi-attribute settings of SymNet, there is a significant performance drop on aPY but a smaller one on SUN without the multi-attribute (correlation) constraint. The reason is the different levels of correlations in aPY and SUN, where aPY is much stronger, as revealed in Fig. 5. It is also noticed that training a SymNet without \mathcal{L}_{sym} , or $\mathcal{L}_{sym} \& \mathcal{L}_{tri}$ leads to a significant degradation for aPY, but a relatively minor drop for the other three datasets.

1-e) Besides, in training, the *utilization order of losses* also makes a difference. In generalized CZSL, if the model is first trained with $\mathcal{L}_{cls} \& \mathcal{L}_{tri}$ and then finetuned with all four losses, the performance is very close to our best model (Top-1 AUC drops from 5.40 to 5.35). However, if the model is first trained with $\mathcal{L}_{axiom} \& \mathcal{L}_{sym}$ and then finetuned with all losses, the Top-1 AUC considerably drops from 5.4 to 4.8 (relatively 11.1%). This phenomenon accords with the above analysis that $\mathcal{L}_{cls}, \mathcal{L}_{tri}$ keeps the basic semantics of embeddings and $\mathcal{L}_{axiom}, \mathcal{L}_{sym}$ further guarantee the rationality of attribute transformations. The subsequent transformations cannot be well learned without the well pre-positioned object embeddings and semantic relationships.

(2) **Attention:** we conduct an ablation study on the attention module by removing the attention and only keep the MLPs. The removal will degrade results on all benchmarks. The model can learn the positions to modify on the object embeddings when operating attribute transformations with attention. We also evaluate different type of attention designs, i.e., using activation function $\tanh(\cdot)$ to convert the attentions into range $(-1, 1)$, but it performs worse than the Sigmoid function. The reason maybe the range of activation function (Sigmoid is $(0, 1)$) and the results show Sigmoid activation is more suitable for the training.

(3) **Distance Metrics:** SymNet with other distance metrics including L_1 and cosine distances are evaluated, i.e., replacing all the distance computation in losses and RMD. They all perform much worse than L_2 . With cosine distance, the accuracy severely drops since cosine distance only measures the angle between embeddings and may not be

enough for complex attribute transformations. Moreover, training SymNet with L_1 is more difficult to converge, thus performing poorly. With the same hyper-parameters, the model with L_1 converges at 4,000 epochs, much slower than the model with L_2 (320 epochs). There are more spikes on the L_1 loss curve, indicating its training instability.

For more please refer to the supplementary: comparisons with AttrOperator [2] (Suppl Sec. 4.2) and COMP [3] (Suppl Sec. 4.3), application details (Suppl Sec. 4.4), the relationship between dataset and performance (Suppl Sec. 5).

5 CONCLUSION

In this work, we study the symmetry property of the attribute-object compositions, which reveals profound principles in composition transformations. To an object, if we add an attribute that it already has, or erase one it does not have, it would keep. We construct a framework inspired by group theory to couple and decouple attribute-object compositions to learn symmetry and use group axioms and symmetry as the learning objectives. Moreover, we explore the attribute correlation to improve attribute recognition with the extended learning objectives with multi-attribute constraints for a multi-attribute scenario. On attribute learning and CZSL tasks, our method achieves state-of-the-art performances. In the future, we consider to study the transformation with varying degrees, e.g., not-, half-, and totally-peeled and apply SymNet to GAN.

ACKNOWLEDGMENT

This work is supported in part by the National Key R&D Program of China, No.2017YFA0700800, National Natural Science Foundation of China under Grants 61772332, and Baidu Scholarship.

REFERENCES

- [1] Ishan Misra, Abhinav Gupta, and Martial Hebert. From red wine to red tomato: Composition with context. In *CVPR*, 2017.
- [2] Tushar Nagarajan and Kristen Grauman. Attributes as operators: factorizing unseen attribute-object compositions. In *ECCV*, 2018.
- [3] Pavel Tokmakov, Yu-Xiong Wang, and Martial Hebert. Learning compositional representations for few-shot recognition. In *ICCV*, 2019.
- [4] Zhixiong Nan, Yang Liu, Nanning Zheng, and Song-Chun Zhu. Recognizing unseen attribute-object pair with generative model. In *AAAI*, 2019.
- [5] Xin Wang, Fisher Yu, Ruth Wang, Trevor Darrell, and Joseph E Gonzalez. Tafe-net: Task-aware feature embeddings for low shot learning. In *CVPR*, 2019.
- [6] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *arXiv preprint arXiv:1301.3781*, 2013.
- [7] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. In *JMLR*, 2008.
- [8] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.
- [9] Yikang Li, Wanli Ouyang, Bolei Zhou, Kun Wang, and Xiaogang Wang. Scene graph generation from objects, phrases and region captions. In *ICCV*, 2017.
- [10] Phillip Isola, Joseph J Lim, and Edward H Adelson. Discovering states and transformations in image collections. In *CVPR*, 2015.
- [11] Aron Yu and Kristen Grauman. Semantic jitter: Dense supervision for visual comparisons via synthetic images. In *ICCV*, 2017.
- [12] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.

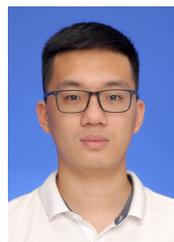
- [13] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *CVPR*, 2009.
- [14] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. In *TPAMI*, 2018.
- [15] Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. In *CVPR*, 2009.
- [16] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [19] Devi Parikh and Kristen Grauman. Relative attributes. In *ICCV*, 2011.
- [20] Yubin Deng, Ping Luo, Chen Change Loy, and Xiaoou Tang. Pedestrian attribute recognition at far distance. In *ACMMM*, 2014.
- [21] Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Describing people: A poselet-based approach to attribute classification. In *ICCV*, 2011.
- [22] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *ICCV*, 2015.
- [23] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. In *arXiv preprint arXiv:1212.0402*, 2012.
- [24] Genevieve Patterson and James Hays. Coco attributes: Attributes for people, animals, and objects. In *ECCV*, 2016.
- [25] Chao-Yeh Chen and Kristen Grauman. Inferring analogous attributes. In *CVPR*, 2014.
- [26] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. In *IJCV*, 2016.
- [27] Sung Ju Hwang, Fei Sha, and Kristen Grauman. Sharing features between objects and their attributes. In *CVPR*, 2011.
- [28] Dhruv Mahajan, Sundararajan Sellamanickam, and Vinod Nair. A joint learning framework for attribute models and object descriptions. In *ICCV*, 2011.
- [29] Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona. Caltech-ucsd birds 200. 2010.
- [30] Yong-Lu Li, Siyuan Zhou, Xijie Huang, Liang Xu, Ze Ma, Hao-Shu Fang, Yanfeng Wang, and Cewu Lu. Transferable interactiveness knowledge for human-object interaction detection. In *CVPR*, 2019.
- [31] Yong-Lu Li, Liang Xu, Xinpeng Liu, Xijie Huang, Yue Xu, Shiyi Wang, Hao-Shu Fang, Ze Ma, Mingyang Chen, and Cewu Lu. PaStaNet: Toward Human Activity Knowledge Engine. In *CVPR*, 2020.
- [32] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. Learning to detect human-object interactions. In *WACV*, 2018.
- [33] Senthil Purushwalkam, Maximilian Nickel, Abhinav Gupta, and Marc’Aurelio Ranzato. Task-driven modular networks for zero-shot compositional learning. In *ICCV*, 2019.
- [34] Wei-Lun Chao, Soravit Changpinyo, Boqing Gong, and Fei Sha. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *ECCV*, 2016.
- [35] Zeynep Akata, Florent Perronnin, Zaid Harchaoui, and Cordelia Schmid. Label-embedding for image classification In *PAMI*, 2015.
- [36] Sun-Wook Choi, Chong Ho Lee, and In Kyu Park. Scene classification via hypergraph-based semantic attributes subnetworks identification. In *ECCV*, 2014.
- [37] Chuang Gan, Tianbao Yang, and Boqing Gong. Learning attributes equals multi-source domain generalization. In *CVPR*, 2016.
- [38] Zhi-Qi Cheng, Xiao Wu, Siyu Huang, Jun-Xiu Li, Alexander G Hauptmann, and Qiang Peng. Learning to transfer: Generalizable attribute learning with multitask neural model search. In *ACMMM*, 2018.
- [39] Kongming Liang, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Unifying visual attribute learning with object recognition in a multiplicative framework. In *PAMI*, 2018.
- [40] Yongxi Lu, Abhishek Kumar, Shuangfei Zhai, Yu Cheng, Tara Javidi, and Rogerio Feris. Fully-adaptive feature sharing in multi-task networks with applications in person attribute classification. In *CVPR*, 2017.
- [41] Emily M Hand, and Rama Chellappa. Attributes for improved attributes: A multi-task network utilizing implicit and explicit relationships for facial attribute classification. In *AAAI*, 2017.
- [42] Yuval Atzmon, Felix Kreuk, Uri Shalit, and Gal Chechik. A causal view of compositional zero-shot recognition. In *arXiv:2006.14610*, 2020.
- [43] Yongqin Xian, Tobias Lorenz, Bernt Schiele, and Zeynep Akata. Feature generating networks for zero-shot learning. In *CVPR*, 2018.
- [44] Jie Yang, Jiarou Fan, Yiru Wang, Yige Wang, Weihao Gan, Lin Liu, and Wei Wu. Hierarchical Feature Embedding for Attribute Recognition. In *arXiv:2005.11576*, 2020.
- [45] Chufeng Tang, Lu Sheng, Zhaoxiang Zhang, and Xiaolin Hu. Improving Pedestrian Attribute Recognition With Weakly-Supervised Multi-Scale Attribute-Specific Localization. In *CVPR*, 2018.
- [46] Yuze Li and Chunling Yang and Yu Chen and Yan Zhang. Un-supervised domain adaptation with structural attribute learning networks. In *Neurocomputing*, 2020.
- [47] Muhammad Ferjad Naeem, Yongqin Xian, Federico Tombari, and Zeynep Akata. Learning Graph Embeddings for Compositional Zero-shot Learning. In *CVarXiv:2102.01987PR*, 2021.
- [48] Yong-Lu Li, Xinpeng Liu, Xiaoqian Wu, Yizhuo Li, and Cewu Lu. HOI analysis: Integrating and decomposing human-object interaction. In *NeurIPS*, 2020.
- [49] Yong-Lu Li, Xinpeng Liu, Han Lu, Shiyi Wang, Junqi Liu, Jiefeng Li, and Cewu Lu. Detailed 2D-3D Joint Representation for Human-Object Interaction. In *CVPR*, 2020.
- [50] Yong-Lu Li, Yue Xu, Xiaohan Mao, and Cewu Lu. Symmetry and Group in Attribute-Object Compositions. In *CVPR*, 2020.



Yong-Lu Li received a Ph.D. degree in computer science and technology from Shanghai Jiao Tong University under the supervision of Prof. Cewu Lu. He has received the Baidu Scholarship, YunFan award, Shanghai Outstanding Graduate, etc. His research interests include computer vision, reasoning, and embodied AI. He has developed HAKE, which is a knowledge-driven perception and reasoning system for human-scene-robot interaction.



Yue Xu received a B.E. degree, and is working toward a master's degree in computer science from Shanghai Jiao Tong University, Shanghai, China. His research interests mainly include computer vision.



Xinyu Xu is an undergraduate student majoring in Computer Science and Engineering, Shanghai Jiao Tong University. His research interests include computer vision and robotics.



Xiaohan Mao is an undergraduate student majoring in Computer Science and Engineering, ACM class, Shanghai Jiao Tong University. His research interests include computer vision.



Cewu Lu is an Associate Professor at Shanghai Jiao Tong University (SJTU). Before joining SJTU, he was a research fellow at Stanford University, working under Prof. Fei-Fei Li and Prof. Leonidas J. Guibas. He was a Research Assistant Professor at Hong Kong University of Science and Technology with Prof. Chi Keung Tang. He got his Ph.D. degree from The Chinese University of Hong Kong, supervised by Prof. Jiaya Jia. His research interests fall mainly in computer vision, deep learning, and robotics.

APPENDIX

This is a PAMI version of our CVPR'20 work SymNet [50].

The specific hyper-parameters in experiments are shown in Tab. 8.

Here, we detail the datasets and baselines adopted in the experiments.

.1 Datasets

aPY [15] consists of two **multi-attribute** datasets, aPascal and aYahoo, both with 64 attributes. aPascal has 6,340 training samples and 6,355 testing samples covering 20 objects. aYahoo contains 2,644 images for testing, covering 12 objects disjoint from aPascal. Following [37], [38], we train our model on aPascal and test on aYahoo.

SUN [12] is a **multi-attribute** dataset, contains 14,340 images with 102 attributes and 717 objects. 10,320 and 4,020 images are used for training and testing respectively.

MIT-States [10] contains 63,440 images covering 245 objects and 115 attributes. Each image is attached with one **single** object-attribute composition label and there are 1,262 possible pairs in total. We follow the setting of [1] and use 1,262 pairs/34,562 images for training and 700 pairs/19,191 images as the test set.

UT-Zappos50K [11] is a fine-grained and **single-attribute** dataset with 50,025 images of shoes annotated with shoe type-material pairs. We follow the setting and split from [2], using 83 object-attribute pairs/24,898 images as the train set and 33 pairs/4,228 images for testing.

.2 Baselines for Single-Attribute Learning and CZSL

Visual Product trains two simple classifiers for attributes and objects independently and fuses the outputs by multiplying their margin probabilities: $P(a, o) = P(a)P(o)$. The classifiers can be either linear SVMs [1] or single layer softmax regression models [2].

LabelEmbed (LE) [1] combines the word vectors [16] of attribute and object and uses 3-layer FCs to transform the pair embedding into a transform matrix. The classification score is the product of transform matrix and visual feature:

- 1) **LabelEmbed Only Regression (LEOR)** [1] changes the target to minimize the Euclidean distance between $\mathcal{T}(e_a, e_b)$ and the weight of pair SVM classifier w_{ab} .
- 2) **LabelEmbed With Regression (LE+R)** [1] combines the losses of LE and LEOR aforementioned.
- 3) **LabelEmbed+** [2] embeds the attribute, object vectors, and image features into a semantic space and also optimizes the input representations during training.

AnalogousAttr [25] trains linear classifiers for seen compositions and uses tensor completion to generalize to the unseen pairs. We report the reproduced results from [2].

Red Wine [1] uses SVM weights as the attribute or object embeddings to replace the word vectors in LabelEmbed.

AttrOperator [2] regards attributes as linear transformations and object word vectors [16] after transformation as pair embeddings. It takes the pair with the closest distance to the image feature as the recognition result. Besides the top-1 accuracy directly reported in [2], we evaluate the top-2, three accuracies with the open-sourced code.

TAFE-Net [5] uses word vectors [6] of attribute-object pair as task embeddings of its meta learner. It generates a binary classifier for each existing composition. We report the results based on VGG-16, which is *better* and more complete than the result based on ResNet-18.

GenModel [4] projects the visual features of images and semantic language embeddings of pairs into a shared latent space. The prediction is given by comparing the distance between visual features and all candidate pair embeddings. **f-CLSWGAN** [43] generates unseen class features via GAN and train a classifier jointly with real and generated features. **TMN** [33] adopts a set of FC-based modules and configure them via a gating function in a task-driven way. It can be generalized to unseen pairs via re-weighting primitive modules.

Causal [42] proposes a causal view of CZSL and learns better representations by disentangling the attribute and object features according to the conditional independence principle. The predictions are given according to the distance between learned features and the attribute/object centers.

.3 Baselines for Multiple-Attribute Learning

ALE [35] embeds objects with category-level attributes. It trains attribute classifiers with an objective to meet correct object embedding. We report the score reproduced by [38].

HAP [36] constructs hyper-graph to learn the correlations of semantic attributes. We report the result by [37].

UDICA/KDICA [37] regularizes the distributional variance to achieve cross-domain attribute generalization. KDICA integrates kernel alignment for a unified optimization. The score on SUN [12] is reproduced by [38].

Dataset	MIT-States [10]	MIT-States (generalized) [10]	UT-Zappos [11]	UT-Zappos (generalized) [11]	aPY [15]	SUN [12]
Learning rate	5e-4	3e-4	1e-4	1e-3	3e-3	5e-3
Batch size	512	512	256	512	128	128
Epoch	320	1000	600	290	177	95
λ_1	5e-2	2e-2	1e-2	2e-2	5e-2	8e-3
λ_2	1e-2	2e-2	3e-2	1e-2	5e-3	1e-3
λ_3	1	1	1	1	1	1
λ_4	1e-2	1e-2	5e-1	1e-2	5e-2	3e-1
λ_5	3e-2	1	5e-1	1	1	5e-2
λ_6	/	/	/	/	5e-2	6e-2
λ_7	/	/	/	/	1	6e-1
Triplet margin	0.5	0.3	0.5	0.5	0.5	0.5

TABLE 8: Hyper-parameters on four benchmarks.

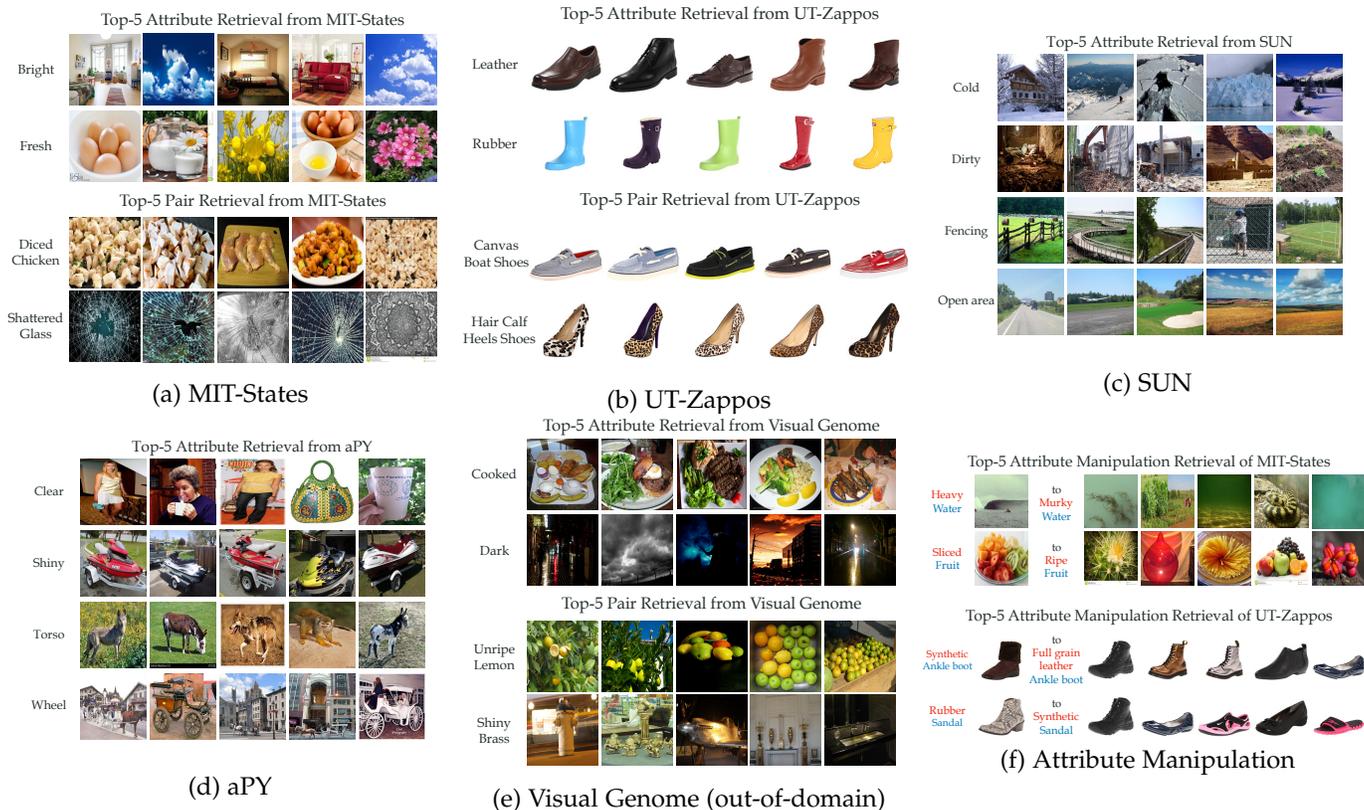


Fig. 9: Image retrieval results. On MIT-States [10], UT-Zappos [11], SUN [12], aPY [15], and Visual Genome [26], the first row shows retrievals of attributes and the second is retrievals of unseen attribute-object pairs. Specially, the retrieval model in out-of-domain [2] mode is not trained on Visual Genome. At last, we also show the retrievals after the attribute manipulation via SymNet.

UMF [39] projects both image features and category labels to a common latent space. Then it makes element-wise multiplication and predicts attributes. We report the score reproduced by [38].

AMF [41] designs a multi-task deep neural network for multiple attribute prediction. It uses an auxiliary network to explore attribute relations further. We report the score reproduced by [38].

FMT [40] automatically designs a neural network that greedily makes branch and task-grouping decisions in each layer. We report the score reproduced by [38].

GALM [38] applies a tree-structured model. Its root node is shared for all attributes, but leave nodes are independently and automatically searched.

To qualitatively evaluate our method, we further report the image retrieval results of SymNet. We follow the settings

of [2]: 1) **In-domain attributes or unseen compositions**: we train SymNet on MIT-States [10], UT-Zappos [11], SUN [12] and aPY [15] respectively, and query the attributes or unseen pairs upon the test set of each dataset. The results are displayed in Fig. 9 (a, b, c, d). 2) **Out-of-domain retrieval**: with SymNet *only trained on MIT-States*, we conduct retrieval on the large-scale Visual Genome [26] with over 100K images, which is non-overlapping with the train set of MIT-States. The results are shown in Fig. 9(e).

Our model is capable of recognizing the images with queried attributes and pairs in most cases. When querying an attribute, it accurately retrieves images across various objects, e.g. for MIT-States, the top-5 retrievals of attribute *fresh* vary among *fresh-egg*, *fresh-milk* and *fresh-flower*, suggesting that our model has well exploited the contextuality and compositionality of at-

	AUC Top-1	AUC Top-2	AUC Top-3	Seen Acc.	Unseen Acc.	H-Mean
SymNet	5.4	11.6	16.6	30.4	25.8	17.6
SymNet w/o \mathcal{L}_{sym}	4.0	9.2	14.0	24.2	24.4	15.4
SymNet w/o \mathcal{L}_{action}	4.0	9.4	14.1	25.3	23.4	15.3
SymNet w/o \mathcal{L}_{inv}	4.3	9.6	14.4	26.0	24.6	15.8
SymNet w/o \mathcal{L}_{com}	4.4	9.6	14.5	26.5	24.8	15.8
SymNet w/o \mathcal{L}_{cls}	4.2	9.9	14.8	25.9	24.3	15.5
SymNet w/o \mathcal{L}_{cls}	1.6	4.7	8.4	13.3	19.9	9.6
SymNet w/o \mathcal{L}_{cls}	3.9	9.7	14.7	25.8	23.5	15.6
SymNet w/o \mathcal{L}_{sym} & \mathcal{L}_{cls}	3.9	9.0	14.1	24.8	24.2	15.1
SymNet w/o \mathcal{L}_{cls} & \mathcal{L}_{cls}	1.6	4.7	8.2	12.7	19.3	9.5
SymNet w/o \mathcal{L}_{sym} & \mathcal{L}_{cls}	1.7	4.7	8.1	13.1	19.5	9.9
SymNet only \mathcal{L}_{sym}	1.6	4.7	8.4	12.4	20.2	9.4
SymNet only \mathcal{L}_{cls} , \mathcal{L}_{cls}	3.7	9.0	13.5	25.0	23.2	15.0
SymNet w/o attention	3.9	9.1	13.8	23.2	24.8	15.1
SymNet tanh attention	3.8	8.7	13.1	24.5	23.6	14.7
SymNet L_1 dist.	2.8	6.3	9.6	22.9	17.8	12.9
SymNet Cos dist.	1.5	4.4	8.0	12.0	20.4	9.1

TABLE 9: Results of ablation studies in generalized CZSL setting [33] on MIT-States [10] validation set.

	MIT-States			UT-Zappos		
	Top-1	Top-2	Top-3	Top-1	Top-2	Top-3
SymNet + SymNet losses (ours)	19.9	28.2	33.8	52.1	67.8	76.0
SymNet + AttrOperator losses	17.0	26.0	31.7	49.8	66.3	73.8
SymNet + Linear operator	16.5	25.5	31.3	49.6	66.2	73.9
AttrOperator + AttrOperator losses	14.2	19.6	25.1	46.2	56.6	69.2
AttrOperator + SymNet losses	14.4	19.9	25.7	46.5	56.7	70.1

TABLE 10: Cross comparison of SymNet and AttrOperator [2].

tributes. In out-of-domain retrieval, SymNet also shows its robustness. Though it has never seen the images in Visual Genome [26], SymNet generalizes well on the target domain and returns correct retrievals, e.g. `dark objects` and `unripe lemon`.

We further report the image retrieval results after attribute manipulation. We first train SymNet on MIT-States [10] or UT-Zappos [11], then use trained CoN and DeCoN to manipulate the image embeddings. For an image with pair label (a, o) , we remove the attribute a with DeCoN and add an attribute b with CoN. Then we retrieve the top-5 nearest neighbors of the manipulated embeddings. This task is much more difficult than the normal attribute-object retrieval [1], [2], [33] because of the complex semantic manipulation and recognition. The results are shown in Fig. 9(f), where the images on the left are original ones and the right ones are the nearest neighbors after manipulation. SymNet is capable of retrieving a certain number of correct samples among the top-5 nearest neighbors, especially in a fine-grained dataset like UT-Zappos [11], suggesting that our model has well exploited the learned symmetry in attribute transformation and learned the contextuality and compositionality of attributes.

.4 Supplementary Ablation Studies

We further conduct more ablation studies under different metrics on multi-dataset to fully evaluate our methods. Precisely, we follow the metric of TMN [33] to conduct the ablations similar to the Tab.7 of the main text. The results are listed in supplementary Tab. 9, which can further verify the effectiveness of the proposed components in SymNet. For more analyses, please refer to Sec. 4.8 of the main text.

.5 Comparison to AttrOperator

Here, we compare SymNet and AttrOperator [2] since they all see attribute change as transformation. However, the design of SymNet is quite different from AttrOperator [2] in many aspects.

First, the operators in AttrOperator are applied on *object Glove* [16] embeddings to compose the anchor representations of compositions for classification, while our CoN and DeCoN are applied on *image representations* and we use relative moving distance for classification. Though both methods regard attributes as manipulations in latent space, the usage of manipulators, the definition of the latent space, the classification paradigm, and the corresponding constraint design are quite different.

Second, as a consequence of different model designs, our losses are naturally based on group axioms and Symmetry property, while the losses of AttrOperator are designed to conform with the linguistic meaning of attributes [2]. Among these losses, \mathcal{L}_{cls} , \mathcal{L}_{inv} , \mathcal{L}_{com} work similar and the others are different. To evaluate the loss difference and effectiveness, we conduct an ablation study of training SymNet model with AttrOperator losses and training AttrOperator model with SymNet losses. The results are shown in Tab. 10. Comparing to the original models, SymNet with AttrOperator losses drops **2.9%** top-1 accuracy, while AttrOperator with SymNet losses instead achieves a slight improvement, indicating our axiom and symmetry-based losses are more complete and especially suitable to SymNet.

Third, we conduct an ablation study about representation power of *linear* model and CoN/DeCoN, as show in Tab. 10. We implement CoN and DeCoN in SymNet as multiple independent linear matrices following [2], which leads to **3.4%** accuracy drop. The experiments show that our CoN and DeCoN are more expressive.

.6 Orthogonal Constraint

Method	aPY	SUN
SymNet (multi)	86.1	88.4
SymNet (multi) + orthogonal	85.9	88.4
SymNet (single)	82.2	88.1
SymNet (single) + orthogonal	84.4	88.3

TABLE 11: Results of ablation studies of orthogonality constraint [3].

Different from the constraints in SymNet, COMP [3] proposed a more straightforward method on multi-attribute learning via orthogonal constraint. To compare their effectiveness, we conduct an extra ablation study on regularizing the orthogonality of **attribute attentions** generated by CoN and DeCoN. Results are reported in Tab. 11.

First, we directly apply the orthogonal loss into the **multi-attribute** setting of SymNet (both the orthogonal loss and our multi-attribute correlation loss are used). It has no noticeable effect and leads to a small degrade on aPY [15]. The reason is that, in the multi-attribute setting of SymNet, we treat each attribute pair concerning their prior attribute correlation in Eq. 10 while orthogonality constraint forces each attribute pair to be orthogonal fairly. Thus, in this scenario, it may not be appropriate to replace the customized attribute relationships with constant orthogonality, e.g., *metal* and *shiny* are strongly correlated attributes that affect objects similarly in transformation and should have similar attribute attentions. Thus, attribute correlation is sound side knowledge in multi-attribute learning.

Second, we apply orthogonal loss into the **single-attribute** setting of SymNet (without the multi-attribute correlation loss). It brings 2.2% and 0.2% performance improvements on aPY [15] and SUN [12] respectively. In this scenario, original attribute attentions are challenging to be distinguished without any regularization. Therefore, orthogonality constraint avoids the dense clusters of attribute attentions and helps SymNet distinguish the different attributes’ roles in coupling and decoupling. Thus, it is a good way to classify multiple attributes without attribute correlation. However, attribute correlation on aPY is very strong, as revealed in Fig. 5 (main text), and there is still a performance gap (1.7%, between 86.1% and 84.4%) between orthogonal constraint and multi-attribute SymNet constraints.

In conclusion, without the attribute correlation, orthogonality constraint is a good alternative way to regularize multiple attributes for SymNet. However, it does not achieve the effect of attribute correlation.

7 Additional Details of Few-Shot Learning

To further verify that whether SymNet can help few-shot recognition learn better representations, we conduct two few-shot recognition experiments on CUB-200-2011 [29] and SUN397 [12] datasets with the same experimental protocol to COMP [3].

CUB-200-2011 dataset [29] is proposed for finer-grained classification of birds. It contains 11,788 images of 200 categories and is split into train and test sets evenly. CUB-200-2011 originally collects 307 category-level attribute annotations. In COMP [3], 130 filtered attributes are utilized since some rare attributes hurt the knowledge transfer in few-shot recognition [3]. Furthermore, the dataset is randomly split into 100 bases and 100 novel categories. We follow their setting here. SUN397 [12] is a scene recognition dataset containing 108,754 images from 397 balanced categories. We use category-level labels of 89 attributes aggregated and filtered by COMP [3]. Moreover, the dataset is randomly split into 197 bases and 200 novel categories.

SymNet is a feature extractor across object categories guided by the group theory constraints. Its attribute transformation feature is rich in attribute semantics and can be used to strengthen the original object representation. Thus, given SymNet trained on two datasets respectively, we can use the concatenation of $f_o, \hat{f}_o^+, \hat{f}_o^-$ as the enhanced attribute representation to the downstream few-shot classification, where f_o is the original feature from ResNet backbone and \hat{f}_o^+, \hat{f}_o^- are the average CoN and DecoN transformed features over all attributes respectively.

In the implementation, we build SymNet on the pre-trained models from COMP [3] as backbones. We run in the two settings from COMP [3], about whether to use data augmentation in training classifier. They are marked as **COMP** and **COMP w/ data aug** respectively in tables in the main text. The default augmentation (indicated with $w/+comp$) in COMP [3] includes *Random Resized Crop*, *Image Jitter*, *Random Horizontal Flip*. We use the base classes in training and set $\lambda_1 = 5e - 2, \lambda_2 = 1e - 3, \lambda_3 = 1, \lambda_4 = 10, \lambda_5 = 1, \lambda_6 = 5e - 2, \lambda_7 = 1$ and learning rate $1e - 3$ on both two datasets. We train SymNet for 1,000 epochs

with batch size 128 on CUB-200-2011 [29], but 500 epochs with batch size 256 on SUN397 [12] since it is much larger. After that, based on the concatenated features ($f_o, \hat{f}_o^+, \hat{f}_o^-$) from SymNet, we train a cosine classifier with learning rate $1e-1$ and batch size 1,000 in 1,2,5-shot settings. It runs for 100 iterations on CUB-200-2011 [29] and 200 iterations on SUN397 [12].

The results on CUB-200-2011 [29] and SUN397 [12] are listed in Tab. 6 (main text) and Tab. 12. Enhanced with SymNet, the COMP [3] baseline gains stable performance improvements on the novel and all categories.

Method	Novel			All		
	1-shot	2-shot	5-shot	1-shot	2-shot	5-shot
COMP [3]	43.4	54.5	65.9	54.9	60.4	66.3
COMP - SymNet	43.6	54.9	66.1	55.5	61.3	67.2

TABLE 12: Supplementary results of few-shot recognition on SUN397 [12].

Analysis of Single-attribute Dataset. Comparatively, accuracy on MIT-States [10] is much lower than UT-Zappos [11] as MIT-States has many more objects and attribute categories and suffers from noisy samples and data insufficiency. Besides, the synonyms and near-synonyms in attributes significantly affect the results. For example, SymNet recognizes 20.4% samples with attribute *ancient* as *old*, while the visual properties of these two attributes can barely be distinguished. These results are correct from the human perspective but mistaken according to the benchmark. To explore this phenomenon on MIT-States, we manually select 13 sets of near-synonyms from MIT-States¹, which are chosen according to the similarity in both linguistic meanings and visual patterns. We then regard the attributes within each set as equal, i.e., predicting the near-synonym is also considered correct. On this new benchmark, our model achieves a 3.03% improvement on attribute accuracy and a 0.66% improvement on CZSL accuracy. We also apply this strategy to AttrOperator [2], obtain an improvement of 2.25% on attribute recognition and 0.28% on CZSL recognition. Comparing to AttrOperator, our model suffers more from the synonym problem.

Analysis of Multi-attribute Dataset. As shown in Fig. 5 of the main text, aPY [15] have stronger overall attribute correlations than SUN [12]. Thus, the ablation study verifies that the contribution of the correlation information in learning is more significant on aPY than SUN. Since images in aPY may have multiple objects, conventional methods that directly crop and resize each instance suffer from noisy boxes and data insufficiency. Statistically, each instance in aPY occupies 25.1% of the image. And instances in the same image are usually of the same category. Therefore, learning from the context may contribute to a better object representation. For aPY, our model with feature extracted by ROI-pooling [8] performs much better than the direct cropping.

1. {cracked, shattered, splintered}; {chipped, cut}; {dirty, grimy}; {eroded, weathered}; {huge, large}; {melted, molten}; {ancient, old}; {crushed, pureed, mashed}; {ripped, torn}; {crinkled, crumpled, ruffled, wrinkled}; {small, tiny; damp, wet}