# TransZero++: Cross Attribute-Guided Transformer for Zero-Shot Learning

Shiming Chen, Ziming Hong, Wenjin Hou, Guo-Sen Xie, Yibing Song, Jian Zhao,
Xinge You, *Senior Member, IEEE*, Shuicheng Yan, *Fellow, IEEE*, and Ling Shao, *Fellow, IEEE*

**Abstract**—Zero-shot learning (ZSL) tackles the novel class recognition problem by transferring semantic knowledge from seen classes to unseen ones. Semantic knowledge is typically represented by attribute descriptions shared between different classes, which act as strong priors for localizing object attributes that represent discriminative region features, enabling significant and sufficient visual-semantic interaction for advancing ZSL. Existing attention-based models have struggled to learn inferior region features in a single image by solely using unidirectional attention, which ignore the transferable and discriminative attribute localization of visual features for representing the key semantic knowledge for effective knowledge transfer in ZSL. In this paper, we propose a cross attribute-guided Transformer network, termed TransZero++, to refine visual features and learn accurate attribute localization for key semantic knowledge representations in ZSL. Specifically, TransZero++ employs an attribute→visual Transformer sub-net (AVT) and a visual→attribute Transformer sub-net (VAT) to learn attribute-based visual features and visual-based attribute features, respectively. By further introducing feature-level and prediction-level semantical collaborative losses, the two attribute-guided transformers teach each other to learn semantic-augmented visual embeddings for key semantic knowledge representations via semantical collaborative learning. Finally, the semantic-augmented visual embeddings learned by AVT and VAT are fused to conduct desirable visual-semantic interaction cooperated with class semantic vectors for ZSL classification. Extensive experiments show that TransZero++ achieves the new state-of-the-art results on three golden ZSL benchmarks and on the large-scale ImageNet dataset. The project website is available at: https://shiming-chen.github.io/TransZero-pp/TransZero-pp.html.

**Index Terms**—Zero-Shot Learning; Transformer; Attribute Localization; Semantic-Augmented Visual Embedding; Semantical Collaborative Learning.

✦

## 1 INTRODUCTION AND MOTIVATION

HUMAN beings are capable of learning novel concepts based on prior experience without seeing them in advance. For example, given the clues that zebras appear like horses yet with black-and-white stripes of tigers, one can quickly recognize a zebra if he/she has seen horses and tigers before. Nevertheless, unlike humans, supervised machine learning models can only classify samples belonging to the classes that have already appeared during the training phase, and they are not capable of handling samples from previously unseen categories. Motivated by this challenge, zero-shot learning (ZSL) was proposed to recognize new classes by exploiting the intrinsic semantic relatedness during learning [1], [2], [3], [4], [5], [6]. Since ZSL is a foundational method of artificial intelligence, it is commonly used in tasks with wide real-world applications, *e.g.*, image classification [7], [8], image retrieval [9], [10], semantic segmentation [11] and object detection [12]. Particularly, the core idea of ZSL is to learn discriminative and transferable visual features for conducting effective visual-semantic interactions based on the semantic information (*e.g.*, attribute descriptions [4], sentence embeddings [13], and DNA [14]), which are shared between the seen and unseen classes

*S. Chen, Z. Hong, W. Hou and X. You are with the School of Electronic Information and Communication, Huazhong University of Science and Technology, Wuhan 430074, China. (Corresponding author: Xinge You. e-mail: youxg@hust.edu.cn)*
*G.-S. Xie is with the Nanjing University of Science and Technology, China.*
*Y. Song is with AI³ Institute, Fudan University, Shanghai, China.*
*J. Zhao is with the Institute of North Electronic Equipment, Beijing, China, and the Peng Cheng Laboratory, Shenzhen, China.*
*S. Yan is with Sea AI Lab (SAIL), Singapore.*
*L. Shao is with Terminus Group, China.*

employed to support the knowledge transfer. At present, most existing ZSL methods bases on attribute descriptions. According to the different ranges of the label space during testing, ZSL methods can be categorized into conventional ZSL (CZSL), which aims to predict unseen classes, and generalized ZSL (GZSL), which can predict both seen and unseen classes [15]. Moreover, ZSL can also be classified as inductive ZSL [16], [17], which only utilizes the labeled seen data, and transductive ZSL [18], [19], assuming that unlabeled unseen data are available [15]. Inductive ZSL is more reasonable and challenging, we are thus focused on the inductive ZSL setting in this paper. ZSL is typically denoted as zero-shot image classification or object recognition [3], [4], which is different to zero-shot text classification [20] or zero-shot transfer of model [21]. We also follow this standard in this paper.

To enable visual-semantic interactions for knowledge transfer from seen to unseen classes, early embedding-based ZSL methods [18], [23], [24] are trying to learn the embedding between seen classes and their class semantic vectors, and then classify unseen classes by nearest neighbor search in the embedding space. However, these embedding-based methods inevitably overfit to seen classes under the GZSL setting (known as the bias problem), since the embedding is only learned by seen class samples. To mitigate this bias problem, many generative ZSL methods have been proposed to synthesize feature samples for unseen classes by leveraging generative models (*e.g.*, variational autoencoders (VAEs) [25], [26], [27], generative adversarial nets (GANs) [8], [16], [28], and generative flows [29]) for data augmentation. Thus the generative ZSL methods can compensate for the lack of training samples of unseen classes and convert ZSL into a supervised classification task.
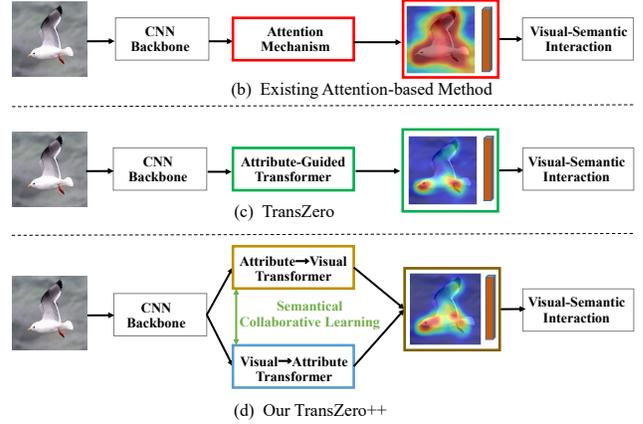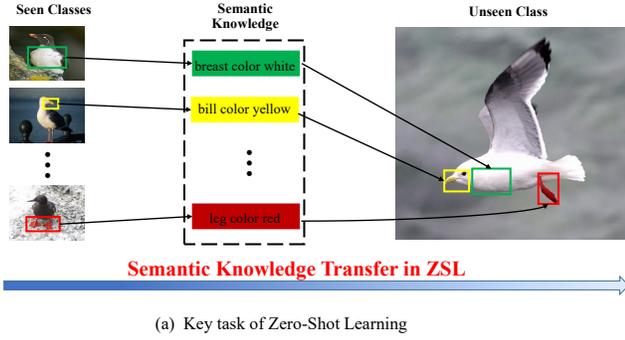
Fig. 1: Motivation illustration. (a) ZSL discovers the discriminative and transferable semantic knowledge to enable efficient knowledge transfer from seen classes to unseen ones. (b) Existing attention-based ZSL methods simply learn inferior region embeddings (e.g., the whole bird body) for semantic knowledge representations using unidirectional attention, ignoring the transferable and discriminative attribute localization (e.g., the distinctive bird body parts) of visual features; (c) TransZero [22] employs an attribute-augmented Transformer to reduce the entangled relationships among region features to improve their transferability, and localizes the object attributes to represent discriminative region features as important semantic knowledge. (d) Our TransZero++ takes two cross attribute-augmented Transformers (*i.e.*, attribute→visual Transformer and visual→attribute Transformer) to further fully discover more intrinsic semantic knowledge via semantical collaborative learning, encouraging more desirable and sufficient visual-semantic interaction.

Although these methods have achieved significant improvement, they rely on global (holistic) visual features[1] which are insufficient for capturing the fine-grained attribute information (*e.g.*, "bill color yellow" of *Red Legged Kittiwake* ) for representing semantic knowledge [34], [35], [36]. Because the discriminative and transferable semantic knowledge is usually contained in a few regions corresponding to a few attributes, enabling efficient knowledge transfer from seen classes to unseen ones, as shown in Fig. 1 (a). Thus, the visual feature representations learned by these methods are inferior, resulting in undesirable visual-semantic interactions for knowledge transfer. More recently, few attention-based models [34], [37], [38], [39], [40], [41], [42] have attempted to explore more discriminative region features, as shown in Fig. 1 (b). However, these methods are limited in: i) they directly take the entangled region (grid) features[2] for ZSL classification, which hinders the transferability of visual features from seen to unseen classes; ii) they simply learn region embeddings (*e.g.*, the whole bird body) using unidirectional attention, neglecting the importance of discriminative attribute localization (*e.g.*, the distinctive bird body parts) for key semantic knowledge representations. Thus, properly discovering key semantic knowledge from visual features to enable efficient semantic knowledge transfer in ZSL has become very necessary.

To tackle the above challenges, in this paper, we propose a cross attribute-guided Transformer, termed TransZero++, which discovers the key semantic knowledge for efficient knowledge transfer from seen classes to unseen ones in ZSL via semantical collaborative learning, as shown in Fig. 1 (d). Specifically, TransZero++ consists of two attribute-guided Transformer sub-nets

(*i.e.*, attribute→visual Transformer (AVT) and visual→attribute Transformer (VAT)) that learn attribute-based visual features and visual-based attribute features respectively, which are further mapped into the semantic embedding space using two mapping functions $\mathcal{M}_1$ and $\mathcal{M}_2$ to conduct desirable visual-semantic interaction. In AVT and VAT, we first take a feature augmentation encoder to augment visual features by i) alleviating the cross-dataset bias between ImageNet and ZSL benchmarks, and ii) reducing the entangled relative geometry relationships between different regions for improving the transferability from seen to unseen classes. These augmented visual features will promote the following sequential learning. To learn locality-augmented visual features, we employ an attribute→visual decoder in AVT to localize the image regions most relevant to each attribute in a given image (denoted as attribute-based visual features), under the guidance of semantic attribute information. We also take a visual→attribute decoder to learn visual-based attribute features in VAT. By introducing feature-level and prediction-level semantical collaborative losses further, the two attribute-guided transformers teach each other to further learn two complementary semantic-augmented visual embeddings for key semantic knowledge representations via semantical collaborative learning. Finally, the semantic-augmented visual embeddings cooperated with the semantic vectors are leveraged to conduct desirable visual-semantic interaction for ZSL classification. Extensive experiments show that TransZero++ achieves the new state-of-the-art on three ZSL benchmarks and on the large-scale ImageNet dataset. The qualitative results also demonstrate that TransZero++ refines visual features and accurately localizes attribute regions for semantic-augmented feature representations.

A preliminary version of this work was presented as a conference paper (termed TransZero [22]). As shown in Fig. 1 (c), although TransZero can localize some important attributes for discriminative region feature representations with low confident scores, some other valuable attributes are failed (*e.g.*, "white wing color" of *Red Legged Kittiwake*). In this version, we strengthen the work from four aspects: i) We propose VAT to learn visual-

---

1. These global visual features are directly extract from a CNN Backbone (*e.g.*, ResNet [30]) pre-trained on ImageNet [31] alone, ignoring the cross-dataset bias between ImageNet and ZSL benchmarks (*e.g.*, CUB [32]). Such a bias inevitably results in poor-quality visual features in which not all the dimensions are semantically related to the pre-defined attributes for ZSL tasks [28], [33].

2. These entangled region (grid) features usually include relative geometry relationship priors between different regions [34].

based semantic attribute representations that are complementary to the attribute-based visual features learned by AVT, enabling TransZero++ to fully discover the key semantic knowledge from visual features. ii) We introduce feature-level and prediction-level semantic collaborative losses to encourage AVT and VAT to teach each other to discover more intrinsic semantic knowledge for improving the confidence scores of attribute localization, under the guidance of semantic collaborative learning. iii) Since the learned attribute-based visual features and visual-based attribute features are complementary to each other, we fuse the two semantic-augmented visual embeddings learned by AVT and VAT to conduct desirable visual-semantic interaction for ZSL classification. iv) We conduct substantially more experiments to demonstrate the effectiveness of the proposed framework and verify the contribution of each component. Thus, TransZero [22] is extended to be TransZero++.

The main contributions of this paper are summarized as follows:

- We introduce a novel ZSL method, termed TransZero++, which simultaneously refines the visual features and localizes the object attributes for semantic knowledge representations via semantical collaborative learning. TransZero++ consists of an attribute→visual Transformer sub-net (AVT) and a visual→attribute Transformer sub-net (VAT) that learns attribute-based visual features and visual-based attribute features, respectively, which are complementary to each other.
- We propose a feature augmentation encoder to i) alleviate the cross-dataset bias between ImageNet and ZSL benchmarks, and ii) reduce the entangled relative geometry relationships between different regions to improve the transferability of visual features. They are ignored by existing ZSL methods. This feature augmentation encoder is incorporated into AVT and VAT.
- We introduce feature-level and prediction-level semantic collaborative losses to enable semantical collaborative learning between the AVT and VAT, encouraging TransZero++ to learn semantic-augmented visual embeddings by discovering more intrinsic semantic knowledge between visual and attribute features.
- Extensive experiments demonstrate that TransZero++ achieves the new state-of-the-art on three popular ZSL benchmarks and on the large-scale ImageNet dataset. Compared with the popular attention-based method (*i.e.*, APN [39]), TransZero++ leads to significant improvements of 6.3%/3.2%, 6.0%/4.9% and 4.2%/8.6% in $acc/H$ on CUB [32], SUN [43] and AWA2 [15], respectively.

The rest of this paper is organized as follows. Sec. 2 discusses related works. The proposed TransZero++ is illustrated in Sec. 3. Experimental results and discussions are provided in Sec. 4. Finally, we present a summary in Sec. 5.

## 2 RELATED WORK

In this section, we mainly review three streams of related works: zero-shot learning, Transformer, and collaborative learning.

### 2.1 Zero-Shot Learning

Early embedding-based ZSL methods [18], [24], [44], [44], [45] aim to learn a mapping from the visual domain to semantic domains to transfer semantic knowledge from seen to unseen classes. They usually extract global visual features from pre-trained or end-to-end trainable networks, *e.g.*, ResetNet [30]. Note that end-to-end models achieve better performance than pre-trained ones because they fine-tune the visual features, thus the cross-dataset bias between ImageNet and ZSL benchmarks is alleviated [8], [28].

However, these methods inevitably overfit to seen classes in GZSL since they only learn the model on seen classes [15], [46], [47]. As such, the generative ZSL methods are introduced to tackle this challenge using various generative models (*e.g.*, VAEs [25], [26], [27], [33], GAN [25], [26], [27], and generative flows [29]) to synthesize a number of images or visual features for unseen classes based on the class semantic vector (attribute values manually annotated by humans). Thus, the ZSL task is converted to supervised classification. Arora *et al.* [25] uses a conditional VAE model (SE-ZSL) to synthesize images for unseen classes. Since synthesizing the high dimensional image is not feasible, Xian *et al.* [8], [16] propose f-CLSWGAN and f-VAEGAN to synthesize visual features based on GANs. Different from these generative methods that learn semantic-to-visual mapping as a generator, the common space learning-based ZSL methods are also a special generative model that maps visual and semantic features into a common space simultaneously using VAEs [26], [33], [48].

Although the aforementioned ZSL methods have achieved significant improvements, they still yield relatively undesirable results. This is because they compress holistic visual features to perform global embedding cannot efficiently capture the subtle differences among various fine-grained classes [35]. Furthermore, the holistic visual features are limited to poor transferable from one domain to another domain (*e.g.*, from seen to unseen classes) [49], [50]. More relevant to this work are the recent attention-based ZSL methods [34], [37], [38], [39], [51] that utilize attribute descriptions as guidance to discover the more discriminative region (or part) features. Unfortunately, They simply learn region embeddings (*e.g.*, the whole bird body) neglecting the importance of discriminative attribute localization (*e.g.*, the distinctive bird body parts). Furthermore, the end-to-end attention models are also time-consuming when it comes to fine-tuning the CNN backbone. In contrast, we propose an attribute-guided Transformer to learn the attribute localization for discriminative region feature representations under non end-to-end ZSL model.

### 2.2 Transformer Model

Transformer models [52], [53], [54], [55] have recently achieved excellent performance on a wide range of language and computer vision tasks, *e.g.*, machine translation [56], image recognition [57], video understanding [58], visual question answering [59], etc. Generally, the success of Transformer can be attributed to its self-supervision and self-attention [54]. The self-supervision enables complex models to be trained without the high cost of manual annotation, which in turn allows generalizable representations that encode useful relationships between the entities presented in a given dataset to be learned. The self-attention layers consider the broad context of a given sequence by learning the relationships between the elements in the token set (*e.g.*, words in the language, or patches in an image). Some methods [58], [60], [61], [62] have also shown that the Transformer can better capture the relationship between various modals (*e.g.*, visual features and language) in parallel during training. Motivated by these, we design an attribute-guided Transformer that reduces the relationships among different regions to improve the transferability of visual features and learns
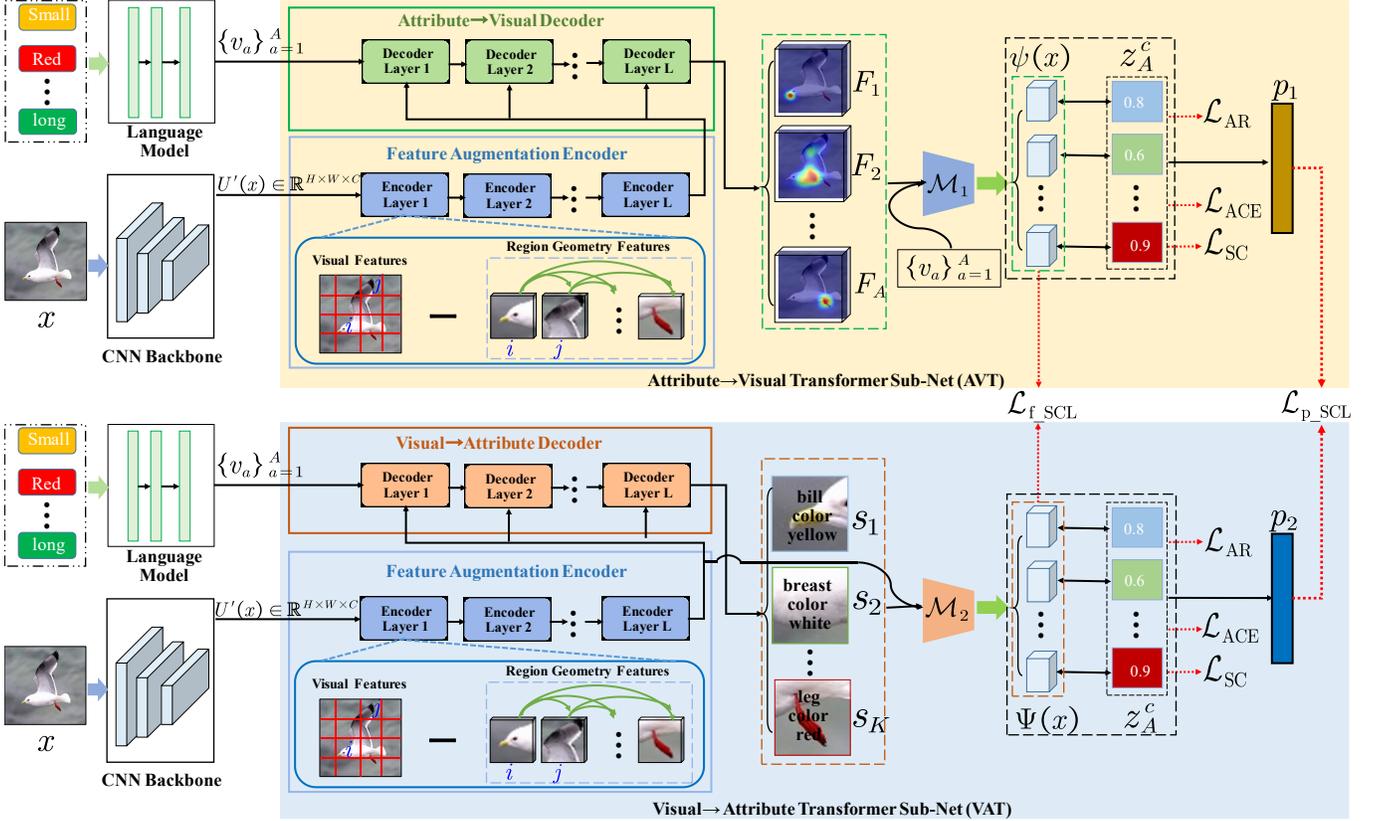
Fig. 2: The architecture of the proposed TransZero++ model. TransZero++ consists of an attribute→visual Transformer sub-net (AVT) and a visual→attribute Transformer sub-net (VAT). AVT includes a feature augmentation encoder that alleviates the cross-dataset bias between ImageNet and ZSL benchmarks and reduces the entangled geometry relationships between different regions for improving the transferability of visual features, and an attribute→visual decoder that localizes object attributes for attribute-based visual feature representations based on the semantic attribute information. Analogously, VAT learns visual-based attribute features using the similar feature augmentation encoder and a visual→attribute decoder. Finally, two mapping functions $\mathcal{M}_1$ and $\mathcal{M}_2$ map the learned attribute-based visual features and visual-based attribute features into semantic embedding space respectively under the guidance of semantical collaborative learning, enabling desirable visual-semantic interaction for ZSL classification.

the attribute localization for representing discriminative region features. In contrast to most of the vision Transformers that learn feature representations on image patches, our TransZero++ learn semantic-augmented visual embeddings on visual features learned by CNN backbone (*e.g.*, ResNet).

## 2.3 Collaborative Learning

Recently, Collaborative Learning [63] has been introduced to learn multiple models jointly for the same task. Teacher-student models to create consistent training supervisions for labeled/unlabeled data using collaborative learning, enabling two-way knowledge transfer from each other. Thus the intrinsic knowledge between different models is distilled for feature representations [64], [65]. Some methods adopt a pool of student models instead of the teacher models by training them with supervision from each other [66], [67]. These works motivate us to design semantical collaborative learning to discover more intrinsic semantic knowledge (*e.g.*, attributes) for semantic-augmented visual embedding representations on the two attribute-guided Transformers. Different from existing collaborative methods that employ multiple similar networks for implicit knowledge distillation, our semantic collaborative learning is based on two attribute-guided Transformers that learn attribute-based visual features and visual-based attribute features respectively for explicit knowledge distillation.

## 3 PROPOSED METHOD

First, we introduce some notations and the problem definition of ZSL. We denote $\mathcal{D}^s = \{(x_i^s, y_i^s)\}$ as training data with $C^s$ seen classes, where $x_i^s \in \mathcal{X}$ refers to the image $i$, and $y_i^s \in \mathcal{Y}^s$ is its corresponding class label. The unseen classes $C^u$ have unlabeled samples $\mathcal{D}^u = \{(x_i^u, y_i^u)\}$, where $x_i^u \in \mathcal{X}$ are the unseen class images, and $y_i^u \in \mathcal{Y}^u$ are the corresponding labels. A set of class semantic vectors of the class $c \in \mathcal{C}^s \cup \mathcal{C}^u = \mathcal{C}$ with $A$ attributes $z^c = [z_1^c, \ldots, z_A^c]^\top = \phi(y)$ (the attribute values are annotated by humans) helps knowledge transfer from seen to unseen classes. According to each word in attribute names, we also take a language model (*i.e.*, GloVe [68]) to learn the semantic attribute features $\mathcal{V}_A = \{v_a\}_{a=1}^A$ of all attributes as auxiliary information ($v_a$ is the $a-th$ attribute feature vector). ZSL aims to predict the class labels $y^u \in \mathcal{Y}^u$ in the CZSL settings and $y \in \mathcal{Y} = \mathcal{Y}^s \cup \mathcal{Y}^u$ in the GZSL setting, where $\mathcal{Y}^s \cap \mathcal{Y}^u = \emptyset$.

In this paper, we propose a cross attribute-guided Transformer network (termed TransZero++) to refine the visual features, localize the object attributes for discriminative region feature representations, and learn semantic-augmented visual embeddings via semantical collaborative learning under a non end-to-end model. This facilitates desirable visual-semantic interaction in ZSL. As illustrated in Fig. 2, our TransZero++ includes an attribute→visual

Transformer sub-net (AVT) and visual→attribute Transformer sub-net (VAT). AVT refines the visual feature using a feature augmentation encoder, and employs an attribute→visual decoder to learn attribute-based visual features, which is further mapped as the semantic-augmented visual embedding $\psi(x)$ in the semantic embedding space using a mapping function $\mathcal{M}_1$. Analogously, VAT uses a similar feature augmentation encoder and a visual→attribute decoder to learn visual-based attribute features, which are further mapped as the semantic-augmented visual embeddings $\Psi(x)$ in the semantic embedding space using another mapping function $\mathcal{M}_2$. Finally, the two semantic-augmented visual embeddings are fused to conduct desirable visual-semantic interaction for ZSL classification based on the class semantic vectors. To encourage TransZero++ to learn semantic-augmented visual embeddings, we introduce feature-level and prediction-level semantical collaborative losses to encourage the two cross AVT and VAT to learn collaboratively and teach each other throughout the training process.

### 3.1 Attribute→Visual Transformer

#### 3.1.1 Feature Augmentation Encoder.

Since the cross-dataset bias between ImageNet and ZSL benchmarks potentially limits the quality of visual feature extraction [28], [33], we first propose a feature augmentation encoder (FAE) to refine the visual features of ZSL benchmarks. In addition, previous ZSL methods simply flatten the grid features $U'(x) \in \mathbb{R}^{H \times W \times C}$ (extracted by a CNN backbone) of a single image into a feature vector, which is further applied to generative models or embedding learning. Unfortunately, such a feature vector implicitly entangles the visual feature representations among various regions in an image, which hinders their transferability from one domain to other domains (e.g., from seen to unseen classes) [39], [59]. Atzmon *et al.* [49] and Chen [50] show that the local visual features are more transferable than the holistic ones. As such, we propose a feature-augmented scaled dot-product attention to further enhance the encoder layer by reducing the relative geometry relationships among the grid features.

Motivated by [69], we first calculate the relative center coordinates $(v_i^{\text{cen}}, t_i^{\text{cen}})$ based on the pair of 2-D relative positions of the $i$-th grid $\{(v_i^{\min}, t_i^{\min}), (v_i^{\max}, t_i^{\max})\}$ during learning the relative geometry features, formulated as:

$$(v_i^{\text{cen}}, t_i^{\text{cen}}) = \left( \frac{v_i^{\min} + v_i^{\max}}{2}, \frac{t_i^{\min} + t_i^{\max}}{2} \right), \quad (1)$$

$$w_i = (v_i^{\max} - v_i^{\min}) + 1, \quad (2)$$

$$h_i = (t_i^{\max} - t_i^{\min}) + 1, \quad (3)$$

where $(v_i^{\min}, t_i^{\min})$ and $(v_i^{\max}, t_i^{\max})$ are the relative position coordinates of the top left corner and bottom right corner of the grid $i$, respectively. Different from [69] that uses 4-D feature vectors for geometry feature representations, we only need to calculate the 2-D geometry feature representations as our grid features are irrelevant to the edges (*i.e.*, width and length) of grid. This is because [69] learn the bounding boxes with various shapes while our grid visual features are shared with same shape.

Different to [69], it attempts to take advantage of the geometric relationship representations for improving the caption performance, by incorporating relative geometry features into the attention weight matrix. In contrast, our FAE aims to remove the relative geometry prior of various image regions from the visual feature representations. Thus, the transferability and discrimination of

visual features are enhanced. Specifically, we construct region geometry features $G_{ij}$ between grid $i$ and grid $j$:

$$G_{ij} = \text{ReLU}\left(w_g^T g_{ij}\right), \quad (4)$$

where

$$g_{ij} = FC\left(r_{ij}\right), \qquad r_{ij} = \begin{pmatrix} \log\left(\frac{|v_i^{\text{cen}} - v_j^{\text{cen}}|}{w_i}\right) \\ \log\left(\frac{|t_i^{\text{cen}} - t_j^{\text{cen}}|}{h_i}\right) \end{pmatrix}, \quad (5)$$

where $r_{ij}$ is the relative geometry relationship between grids $i$ and $j$, $FC$ is a fully connected layer followed by a $ReLU$ activation, and $w_g^T$ is a set of learnable weight parameters.

Finally, we substract the region geometry features from the visual features in the feature-augmented scaled dot-product attention to provide a more accurate attention map, formally defined as:

$$Q^e = U(x)W_q^e, K^e = U(x)W_k^e, V^e = U(x)W_v^e, \quad (6)$$

$$Z_{aug} = \text{softmax}\left(\frac{Q^e K^{e\top}}{\sqrt{d^e}} - G\right)V^e, \quad (7)$$

$$U_{aug}(x) \leftarrow U(x) + Z_{aug}, \quad (8)$$

where $Q$, $K$, $V$ are the query, key, and value matrices, $W_q^e$, $W_k^e$, $W_v^e$ are the learnable matrices of weights, $d^e$ is a scaling factor, and $Z_{aug}$ is the augmented features. $U(x) \in \mathbb{R}^{C \times HW}$ are the packed visual features, which are learned from the flattened features embedded by a fully connected layer followed by a ReLU and a Dropout layer. $U_{aug}(x)$ is the augmented visual features from the feature augmentation encoder, they will promote the following sequential learning. We rewrite $U_{aug}(x)$ as $U_{aug}^{a \rightarrow v}(x)$ and $U_{aug}^{v \rightarrow a}(x)$ in AVT and VAT, respectively.

#### 3.1.2 Attribute→Visual Decoder.

To learn attribute-based visual features, we design attribute→visual decoder to localize the most relevant image regions to the corresponding attributes to extract attribute-based visual features from a given image for each attribute. We can attend to image regions with respect to each attribute, and compare each attribute to the corresponding attended visual region features to determine the importance of each attribute. In the standard Transformer [52], it takes the *self-attention* operator to consider all pairwise relations among the visual feature elements. The operator adjusts every single element by attending them to the others. (as shown in Fig. 3 (Left)). However, our *cross-attention* operator attends to visual features from attribute features, as shown in Fig. 3 (Middle). In the decoding process, the attribute→visual decoder continuously localizes the local visual information under the guidance of semantic attribute features $\mathcal{V}_A$. Thus, our attribute→visual decoder can effectively localize the image regions most relevant to each attribute in a given image. The multi-head cross-attention layer uses the outputs of the encoder $U_{aug}^{a \rightarrow v}(x)$ as keys ($K_t^{a \rightarrow v}$) and values ($V_t^{a \rightarrow v}$) and a set of learnable semantic embeddings $\mathcal{V}_A$ as queries ($Q_t^{a \rightarrow v}$). It is defined as:

$$Q_t^{a \rightarrow v} = \mathcal{V}_A W_{qt}^{a \rightarrow v}, \quad (9)$$

$$K_t^{a \rightarrow v} = U_{aug}^{a \rightarrow v}(x)W_{kt}^{a \rightarrow v}, \quad (10)$$

$$V_t^{a \rightarrow v} = U_{aug}^{a \rightarrow v}(x)W_{vt}^{a \rightarrow v}, \quad (11)$$

$$\text{head}_t = \text{softmax}\left(\frac{Q_t^d K_t^{a \rightarrow v\top}}{\sqrt{\tau}}\right)V_t^{a \rightarrow v}, \quad (12)$$

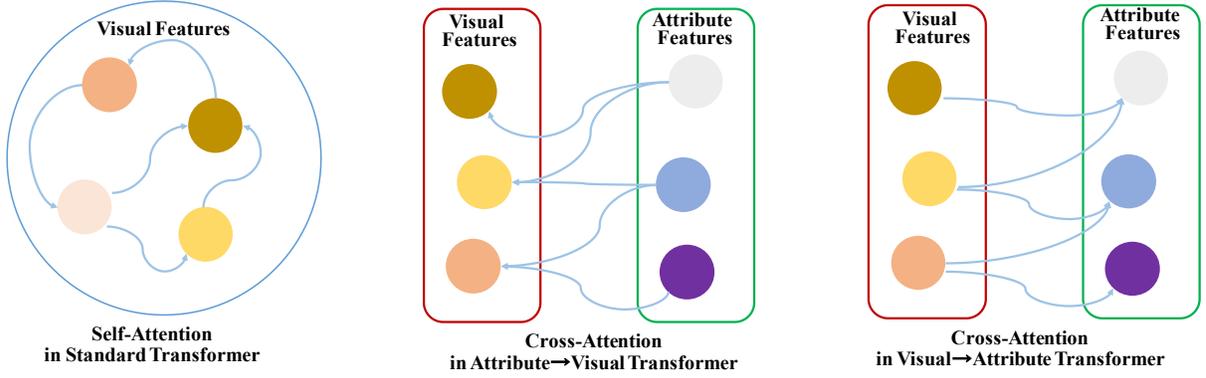$$\hat{F} = \|_{t=1}^{T}(\text{head}_t)W_o^{a \rightarrow v}, \quad (13)$$

Fig. 3: The self-attention operator in standard Transformer, and the cross-attention operator in our attribute→visual Transformer and visual→attribute Transformer.

where $W_{qt}^{a \to v}$, $W_{kt}^{a \to v}$, $W_{vt}^{a \to v}$, and $W_o^{a \to v}$ are the learnable weights, $\tau$ is a scaling factor, and $\|$ is a concatenation operation. Thus, we get a set of attribute-based visual features $\hat{F} = \{\hat{F}_1, \cdots, \hat{F}_A\}$, which captures the visual evidence used to localize the corresponding semantic attributes in the image. Specifically, our AVT will assign a high positive score to the $a$-th attribute if an image has an obvious attribute $v_a$. Otherwise, AVT will assign a low score to the $a$-th attribute. Then, a feed-forward network (FFN) with two linear transformations followed a ReLU activation in between is applied to the attended features $\hat{F}$:

$$F = ReLu\left(\hat{F}W_1^{a \to v} + b_1^{a \to v}\right)W_2^{a \to v} + b_2^{a \to v}, \quad (14)$$

where $W_1^{a \to v}$, $W_2^{a \to v}$, $b_1^{a \to v}$ and $b_2^{a \to v}$ are the weights and biases of the linear layers respectively, and $F = \{F_1, \cdots, F_A\}$ are the final attribute-based visual features that will be mapped into semantic embedding space for desirable visual-semantic interaction using a mapping function ($\mathcal{M}_1$).

### 3.1.3 Visual-Semantic Embedding Mapping

After learning attribute-based visual features that are locality-augmented, we further map them into the semantic embedding space. Based on a mapping function ($\mathcal{M}_1$), we take the semantic attribute vectors $\mathcal{V}_A = \{v_a\}_{a=1}^A$ as support to encourage the mapping to be more accurate. Specifically, $\mathcal{M}_1$ matches the attribute-based visual features $F$ with the semantic attribute information $v_a$, formulated as:

$$\psi(x_i) = \mathcal{M}_1(F) = \mathcal{V}_A^\top W_3^{a \to v} F, \quad (15)$$

where $W_3^{a \to v}$ is an embedding matrix that embeds $F$ into the semantic attribute space. Similar to the class semantic vector $z^c$, $\psi_a(x_i)$ is an attribute score that represents the confidence of having the $a$-th attribute in the image $x_i$. Given a set of semantic attribute vectors $\mathcal{V}_A = \{v_a\}_{a=1}^A$, TransZero++ obtains a mapped semantic-augmented visual embedding $\psi(x_i)$ of a single image $x_i$.

## 3.2 Visual→Attribute Transformer

Likewise, we introduce visual→attribute Transformer (VAT) to attend to attributes with respect to each image region, and thus the visual-based attribute features are learned. They are complementary to the attribute-based visual features, enabling them to calibrate each other to discover more intrinsic semantic knowledge between visual and attribute features. VAT first applies the similar feature augmentation encoder to refine the visual features as $U_{aug}^{v \to a}(x)$, which are further used in visual→attribute decoder of VAT.

### 3.2.1 Visual→Attribute Decoder.

After refining the visual features, we design a visual→attribute decoder to learn visual-based attribute features. Specifically, we take the cross-attention operator to attend the attribute from visual representations, as shown in Fig. 3 (Right). Formally, it is formulated as:

$$Q_t^{v \to a} = U_{aug}^{v \to a}(x)W_{qt}^{v \to a}, \quad (16)$$

$$K_t^{v \to a} = \mathcal{V}_A W_{kt}^{v \to a}, \quad (17)$$

$$V_t^{v \to a} = \mathcal{V}_A W_{vt}^{v \to a}, \quad (18)$$

$$\text{head}_t = \text{softmax}\left(\frac{Q_t^d K_t^{v \to a \top}}{\sqrt{\tau}}\right) V_t^{v \to a}, \quad (19)$$

$$\hat{S} = \|_{t=1}^T (\text{head}_t)W_o^{v \to a}, \quad (20)$$

where $W_{qt}^{v \to a}$, $W_{kt}^{v \to a}$, $W_{vt}^{v \to a}$, and $W_o^{v \to a}$ are the learnable weights, and $\|$ is a concatenation operation. As such, we get a set of visual-based attribute features $\hat{S} = \{\hat{S}_1, \cdots, \hat{S}_K\}$. Intrinsically, $\hat{S}$ is the visual semantic representations corresponding to the $K = H \times W$ visual regions in a single image. Specifically, our VAT will assign a high positive score to the $k$-th visual region with respect to the corresponding attribute, otherwise, VAT will assign a low score. Similar to the AVT, an FFN is utilized to the attended features $\hat{S}$:

$$S = ReLu\left(\hat{S}W_1^{v \to a} + b_1^{v \to a}\right)W_2^{v \to a} + b_2^{v \to a}, \quad (21)$$

where $W_1^{v \to a}$, $W_2^{v \to a}$, $b_1^{v \to a}$ and $b_2^{v \to a}$ are the weights and biases of the linear layers respectively, and $S = \{S_1, \cdots, S_K\}$ are the final visual-based attribute features, which will be mapped into semantic embedding space for significant visual-semantic interaction using a mapping function $\mathcal{M}_2$.

### 3.2.2 Visual-Semantic Embedding Mapping

Once visual-based attribute features are learned, we map them into the semantic embedding space based on a mapping function ($\mathcal{M}_2$). To conduct an effective map, we take the augmented visual features $U_{aug}^{v \to a}(x)$ learned by feature augmentation encoder as support. Thus, $\mathcal{M}_2$ first maps the visual-based attribute features $S$ into $K$ region scores $\bar{S}$, formulated as:

$$\bar{S} = \mathcal{M}_2(S) = U_{aug}^{v \to a}(x)^\top W_3^{v \to a} S, \quad (22)$$

where $W_3^{v \to a}$ is a learnable mapping matrix. Here, $\bar{S}$ is $K$-D, which is not matched with the dimension of class semantic vector

$A$-$D$. Thus, $\mathcal{M}_2$ further embeds $\bar{S}$ into the semantic attribute space with dimension of $A$ based on an attention score $Att = \mathcal{V}_A^\top W_{att} U(x) \in \mathbb{R}^{A \times K}$, where $W_{att}$ is a learnable matrix, written as:

$$\Psi(x_i) = Att \times \bar{S}, \tag{23}$$

Similar to the $\psi(x_i)$, $\Psi_a(x_i)$ is an attribute score that represents the confidence of having the $a$-th attribute in the image $x_i$. As such, TransZero++ obtains a mapped semantic-augmented visual embedding $\Psi(x_i)$ of a single image $x_i$ in VAT.

### 3.3 Model Optimization

To achieve effective optimization for TransZero++, each attribute-guided Transformer sub-net is trained with three supervised losses that have been used in our conference version [22], *i.e.*, the attribute regression loss, attribute-based cross-entropy loss, and self-calibration loss. To enable semantic collaborative learning between the two attribute-guided Transformer sub-nets, *i.e.*, AVT and VAT, we introduce a feature-level semantic collaborative loss and prediction-level semantic collaborative loss, which align each other's visual embeddings and class posterior probabilities respectively.

**Attribute Regression Loss.** To encourage $\mathcal{M}_1$ and $\mathcal{M}_2$ to accurately map visual/attribute features into their corresponding class semantic vectors, we introduce an attribute regression loss to constrain TransZero++. Here, we regard visual-semantic mapping as a regression problem and minimize the mean square error between the embedded attribute score $f(x_i)$ and the corresponding ground truth attribute score $z^c$ of a batch of $n_b$ images $\{x_i^s\}_{i=1}^{n_b}$:

$$\mathcal{L}_{AR} = \frac{1}{n_b} \sum_{i=1}^{n_b} \|f(x_i^s) - z^c\|_2^2. \tag{24}$$

where $f(x_i^s) = \psi(x_i^s)$ for AVT and $f(x_i^s) = \Psi(x_i^s)$ for VAT.

**Attribute-Based Cross-Entropy Loss.** Since the associated visual/attribute embedding is projected near its class semantic vector $z^c$ when an attribute is visually present in an image, we take the attribute-based cross-entropy loss $\mathcal{L}_{ACE}$ to optimize the parameters of the TransZero++, i.e., the dot product between the visual/attribute embedding and each class semantic vector is calculated to produce class logits. This encourages the image/attribute to have the highest compatibility score with its corresponding class semantic vector. Given a batch of $n_b$ training images $\{x_i^s\}_{i=1}^{n_b}$ with their corresponding class semantic vectors $z^c$, $\mathcal{L}_{ACE}$ is defined as:

$$\mathcal{L}_{ACE} = -\frac{1}{n_b} \sum_{i=1}^{n_b} \log \frac{\exp\left(f(x_i^s) \times z^c\right)}{\sum_{\hat{c} \in \mathcal{C}} \exp\left(f(x_i^s) \times z^{\hat{c}}\right)}. \tag{25}$$

**Self-Calibration Loss.** Since $\mathcal{L}_{AR}$ and $\mathcal{L}_{ACE}$ optimize the model only on seen classes, TransZero++ inevitably overfits to these classes [35], [38], [39]. To tackle this challenge, we further employ a self-calibration loss $\mathcal{L}_{SC}$ to explicitly shift some of the prediction probabilities from seen to unseen classes. $\mathcal{L}_{SC}$ is thus formulated as:

$$\mathcal{L}_{SC} = -\frac{1}{n_b} \sum_{i=1}^{n_b} \sum_{c'=1}^{\mathcal{C}^u} \log \frac{\exp\left(f(x_i^s) \times z^{c'} + \mathbb{I}_{[c' \in \mathcal{C}^u]}\right)}{\sum_{\hat{c} \in \mathcal{C}} \exp\left(f(x_i^s) \times z^{\hat{c}} + \mathbb{I}_{[\hat{c} \in \mathcal{C}^u]}\right)}, \tag{26}$$

where $\mathbb{I}_{[c \in \mathcal{C}^u]}$ is an indicator function (*i.e.*, it is 1 when $c \in \mathcal{C}^u$, otherwise -1). Intuitively, $\mathcal{L}_{SC}$ encourages non-zero probabilities to be assigned to the unseen classes during training, which allows

---

**Algorithm 1** The algorithm of TransZero++.

**Input:** The training set $\mathcal{D}^s = \{(x_i^s, y_i^s)\}_{i=1}^{N^s}$, the test set $\mathcal{D}^u = \{(x_i^u, y_i^u)\}_{i=1}^{N^u}$, the pretrained CNN backbone ResNet101, the maximum iteration epoch $\max_{iter}$, loss weights (*i.e.*, $\mathcal{L}_{AR}$, $\mathcal{L}_{SC}$, $\lambda_{VAT}$, $\lambda_{f\_SCL}$ and $\lambda_{p\_SCL}$), the combination coefficients $\alpha$, and hyperparameters (learning rate = 0.001, betas = (0.5, 0.999)) of the Adam optimizer.

**Output:** The predicted label $c^*$ for the test samples.
1: **while** $iter \leq \max_{iter}$ **do** ▷ *Optimization*
2:      Take CNN backbone (*e.g.*, ResNet101 [30]) to extract the visual features $U(x)$ for all image samples.
3:      Take language model (*i.e.*, GloVe [68]) to learn the semantic attribute features $\mathcal{V}_A = \{v_1, \cdots, v_a\}_{a=1}^A$ for each attribute.
4:      Optimize TransZero++ with Eq. 31.
5: **end while**
6: Predict the label $c^*$ of the test samples using Eq. 32. ▷ *Prediction*

---

TransZero++ to produce a large/non-zero probability for the true unseen class when given test samples from unseen classes.

**Semantical Collaborative Loss.** To enable the two attribute-augmented Transformer sub-nets to learn collaboratively and teach each other throughout the training process, we further introduce a feature-level semantical collaborative loss $\mathcal{L}_{f\_SCL}$ and a prediction-level semantical collaborative loss $\mathcal{L}_{p\_SCL}$ for optimization. These two losses are based on $\ell_2$ distance. Note that the $\ell_2$ distance can be replaced with other metrics, *e.g.*, the Kullback Leibler (KL) Divergence or Jensen-Shannon Divergence (JSD).

Specifically, $\mathcal{L}_{f\_SCL}$ uses an $\ell_2$ distance between the semantic-augmented visual embeddings of AVT and VAT (*i.e.*, $\psi(x_i)$ and $\Psi(x_i)$) for test sample $x_i$, formulated as:

$$\mathcal{L}_{f\_SCL} = \frac{1}{n_b} \sum_{i=1}^{n_b} \|\psi(x_i) - \Psi(x_i)\|_2^2. \tag{27}$$

Likewise, $\mathcal{L}_{p\_SCL}$ calculates the $\ell_2$ distance between the predictions of the two attribute-augmented Transformer sub-nets (*i.e.*, $p_1$ and $p_2$), formulated as:

$$\mathcal{L}_{p\_SCL} = \frac{1}{n_b} \sum_{i=1}^{n_b} \|p_1(x_i) - p_2(x_i)\|_2^2, \tag{28}$$

Similar to the TransZero, the AVT and VAT are optimized with the three supervised losses, *i.e.*, $\mathcal{L}_{ACE}$, $\mathcal{L}_{AR}$ and $\mathcal{L}_{SC}$, formulated as:

$$\mathcal{L}_{AVT} = \mathcal{L}_{ACE}^{AVT} + \lambda_{AR}\mathcal{L}_{AR}^{AVT} + \lambda_{SC}\mathcal{L}_{SC}^{AVT}, \tag{29}$$
$$\mathcal{L}_{VAT} = \mathcal{L}_{ACE}^{VAT} + \lambda_{AR}\mathcal{L}_{AR}^{VAT} + \lambda_{SC}\mathcal{L}_{SC}^{VAT}, \tag{30}$$

where $\lambda_{AR}$ and $\lambda_{SC}$ are the loss weights to control the loss $\mathcal{L}_{AR}$ and $\mathcal{L}_{SC}$, respectively, in the AVT and VAT. Finally, we formulate the overall loss function of TransZero++:

$$\begin{aligned}\mathcal{L}_{total} = {} & \mathcal{L}_{AVT} + \lambda_{VAT}\mathcal{L}_{VAT} \\ & + \lambda_{f\_SCL}\mathcal{L}_{f\_SCL} + \lambda_{p\_SCL}\mathcal{L}_{p\_SCL},\end{aligned} \tag{31}$$

where $\lambda_{VAT}$, $\lambda_{f\_SCL}$ and $\lambda_{p\_SCL}$ are the weights to control their corresponding loss terms. To enable the training process of TransZero++ more stable, we set the loss weight to one for $\mathcal{L}_{AVT}$.
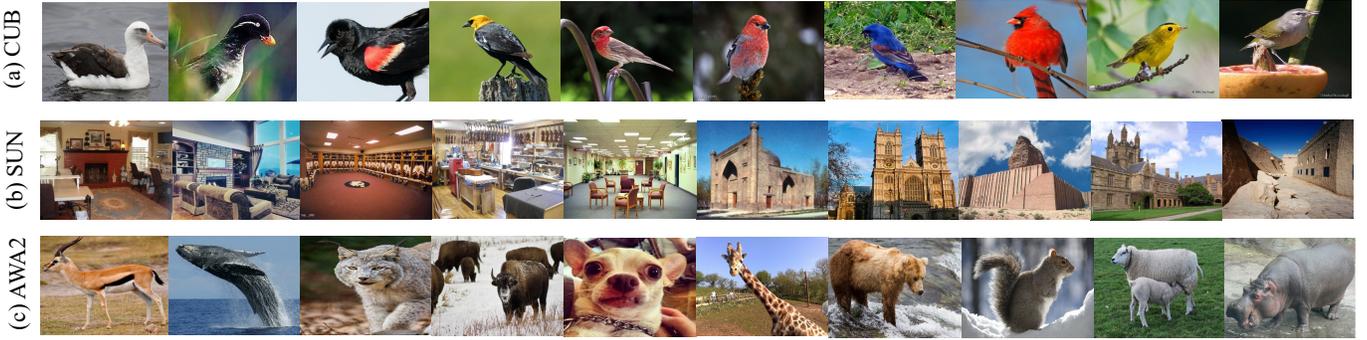
Fig. 4: Some samples on various datasets, including two fine-grained datasets (*i.e.*, CUB and SUN), and one coarse-grained dataset (*i.e.*, AWA2). Each sample is extract from various classes. (Best viewed in color.)

## 3.4 Zero-Shot Prediction

After training TransZero++, We first obtain the visual embeddings of a test instance $x_i$ in the semantic space w.r.t. AVT and VAT, *i.e.*, $\psi(x)$ and $\Psi(x)$. Considering the semantic-augmented visual embeddings learned by AVT and VAT are complementary to each other, we fuse their predictions using a combination coefficients $\alpha$ to predict the test label of $x_i$ with an explicit calibration, formulated as:

$$c^* = \arg \max_{c \in \mathcal{C}^u / \mathcal{C}} (\alpha \psi(x_i) + (1-\alpha)\Psi(x_i))^\top \times z^c + \mathbb{I}_{[c \in \mathcal{C}^u]}. \tag{32}$$

Here, $\mathcal{C}^u/\mathcal{C}$ corresponds to the CZSL/GZSL setting, respectively. The complete procedures (including model training and prediction) for TransZero++ are illustrated by the pseudocode in Algorithm 1.

## 4 EXPERIMENTS

**Dataset.** We conduct extensive experiments on four ZSL benchmarks, including two fine-grained datasets (*i.e.*, CUB [32] and SUN [43]), a coarse-grained dataset (*i.e.*, AWA2 [15]), and a large-scale dataset (*i.e.*, ImageNet [31]). CUB has 11,788 images of 200 bird classes (seen/unseen classes = 150/50) depicted with 312 attributes. SUN includes 14,340 images from 717 scene classes (seen/unseen classes = 645/72) depicted with 102 attributes. AWA2 consists of 37,322 images from 50 animal classes (seen/unseen classes = 40/10) depicted with 85 attributes. We show some samples on these datasets in Fig. 4.

**Evaluation Protocols.** Following [15], we measure the top-1 accuracy both in the CZSL and GZSL settings. In the CZSL setting, we predict the unseen classes to compute the accuracy of test samples, i.e., ***acc***. In the GZSL setting, we calculate the accuracy of the test samples from both the seen classes (denoted as $\boldsymbol{S}$) and unseen classes (denoted as $\boldsymbol{U}$). Meanwhile, their harmonic mean (defined as $\boldsymbol{H = (2 \times S \times U)/(S + U)}$) is also used for evaluation in the GZSL setting.

**Implementation Details.** We use the training splits proposed by [15]. We take a ResNet101 pre-trained on ImageNet as the CNN backbone to extract the visual feature map $U'(x) \in \mathbb{R}^{H \times W \times C}$ ($H$ and $W$ are the height and width of the feature maps, $C$ is the number of channels) without fine-tuning. We use the Adam optimizer with hyperparameters (learning rate = 0.001, betas = (0.5, 0.999)) to optimize our model. The learning rate and batch size are set to 0.0001 and 50 for all datasets, respectively. Following

APN [39], hyperparameters in our model are obtained by grid search on the validation set [15]. Since the training data for the ZSL model is a medium scale which leads to over-fitting with more complex Transformer architectures, the encoder and decoder layers are set to 1 with one attention head both in AVT and VAT. We use PyTorch [89] for the implementation of all experiments. More hyperparameter settings are shown in Sec. 4.6.

## 4.1 Comparison with State-of-the-Art

Our TransZero++ is a non end-to-end and non-generative manner. We compare it with other state-of-the-art methods both in CZSL and GZSL settings, including end-to-end methods (*e.g.*, QFSL [18], SGMA [38], AREN [37], LFGAA [41], APN [39], GEM-ZSL [51]), generative methods (*e.g.*, SE-ZSL [25], f-VAEGAN [8], OCD-CVAE [70], Composer [71], E-PGN [90], TF-VAEGAN [46], IZF [29], SDGZSL [33], GCM-CF [72], FREE [28], HSVA [27], FREE+ESZSL [74]) and non-generative methods (*e.g.*, SP-AEN [76], PQZSL [77], IIR [78], TCN [79], DVBE [80], DAZLE [35]), to demonstrate its effectiveness and advantages.

### 4.1.1 Conventional Zero-Shot Learning

Here, we first compare our TransZero with the state-of-the-art methods in the CZSL setting. As shown in Table 1, our TransZero++ achieves the best accuracy of 78.3% and second-best accuracies of 67.6%/64.6 on CUB and SUN/AWA2, respectively. This shows that TransZero++ effectively discovers the transferable semantic knowledge to represent the locality-augmented features, enabling desirable knowledge transfer for distinguishing various unseen classes. Compared with other attention-based methods (*e.g.*, SGMA [38], AREN [37], APN [39]), TransZero++ gets significant gains over 6.3% and 5.0% on CUB and SUN, respectively. This demonstrates that the attribute localization representations learned by our TransZero++ are more transferable than the region embeddings learned by the existing attention-based methods on fine-grained datasets. Because TransZero++ represents the key transferable semantic knowledge with high confidence, which significantly suppresses the common knowledge between various fine-grained classes. Benefiting from the semantic collaborative learning and the two complementary semantic-augmented embeddings learned by AVT and VAT, TransZero++ further improves the performance over its conference version (TransZero [22]).

To further validate the effectiveness of TransZero++, we compare it with existing methods on large-scale dataset (*i.e.*,

TABLE 1: Results (%) of the state-of-the-art CZSL and GZSL modes on CUB, SUN and AWA2, including end-to-end and non end-to-end methods (generative and non-generative methods). The best and second-best results are marked in **bold** and <u>underline</u>, respectively. The Symbol "–" indicates no results. The symbol "*" denotes attention-based methods. The symbol "†" denotes the methods using calibration. The symbol "‡" denotes the methods using finetuned features.

| Methods | CUB | | | | SUN | | | | AWA2 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CZSL | GZSL | | | CZSL | GZSL | | | CZSL | GZSL | | |
| | acc | U | S | H | acc | U | S | H | acc | U | S | H |
| **End-to-End** | | | | | | | | | | | | |
| QFSL [18] | 58.8 | 33.3 | 48.1 | 39.4 | 56.2 | 30.9 | 18.5 | 23.1 | 63.5 | 52.1 | 72.8 | 60.7 |
| LDF [24] | 67.5 | 26.4 | **81.6** | 39.9 | – | – | – | – | 65.5 | 9.8 | 87.4 | 17.6 |
| SGMA* [38] | 71.0 | 36.7 | 71.3 | 48.5 | – | – | – | – | 68.8 | 37.6 | 87.1 | 52.5 |
| AREN* [37] | 71.8 | 38.9 | 78.7 | 52.1 | 60.6 | 19.0 | 38.8 | 25.5 | 67.9 | 15.6 | <u>92.9</u> | 26.7 |
| LFGAA* [41] | 67.6 | 36.2 | <u>80.9</u> | 50.0 | 61.5 | 18.5 | 40.0 | 25.3 | 68.1 | 27.0 | **93.4** | 41.9 |
| APN*† [39] | 72.0 | 65.3 | 69.3 | 67.2 | 61.6 | 41.9 | 34.0 | 37.6 | 68.4 | 57.1 | 72.4 | 63.9 |
| GEM-ZSL*† [51] | <u>77.8</u> | 64.8 | 77.1 | **70.4** | 62.8 | 38.1 | 35.7 | 36.9 | 67.3 | **64.8** | 77.5 | 70.6 |
| **Non End-to-End** | | | | | | | | | | | | |
| *Generative Methods* | | | | | | | | | | | | |
| SE-ZSL [25] | 59.6 | 41.5 | 53.3 | 46.7 | 63.4 | 40.9 | 30.5 | 34.9 | 69.2 | 58.3 | 68.1 | 62.8 |
| f-CLSWGAN [16] | 57.3 | 43.7 | 57.7 | 49.7 | 60.8 | 42.6 | 36.6 | 39.4 | 68.2 | 57.9 | 61.4 | 59.6 |
| f-VAEGAN [8] | 61.0 | 48.4 | 60.1 | 53.6 | 64.7 | 45.1 | 38.0 | 41.3 | 71.1 | 57.6 | 70.6 | 63.5 |
| OCD-CVAE [70] | – | 44.8 | 59.9 | 51.3 | – | 44.8 | 42.9 | 43.8 | – | 59.5 | 73.4 | 65.7 |
| Composer [71] | 69.4 | 56.4 | 63.8 | 59.9 | 62.6 | **55.1** | 22.0 | 31.4 | 71.5 | 62.1 | 77.3 | 68.8 |
| TF-VAEGAN [46] | 64.9 | 52.8 | 64.7 | 58.1 | 66.0 | 45.6 | 40.7 | 43.0 | 72.2 | 59.8 | 75.1 | 66.6 |
| IZF [29] | 67.1 | 52.7 | 68.0 | 59.4 | **68.4** | <u>52.7</u> | <u>57.0</u> | <u>54.8</u> | **74.5** | 60.6 | 77.5 | 68.0 |
| GCM-CF [72] | – | 61.0 | 59.7 | 60.3 | – | 47.9 | 37.8 | 42.2 | – | 60.4 | 75.1 | 67.0 |
| SDGZSL [33] | 75.5 | 59.9 | 66.4 | 63.0 | – | – | – | – | 72.1 | <u>64.6</u> | 73.6 | 68.8 |
| FREE [28] | – | 55.7 | 59.9 | 57.7 | – | 47.4 | 37.2 | 41.7 | – | 60.4 | 75.4 | 67.1 |
| HSVA [27] | 62.8 | 52.7 | 58.3 | 55.3 | 63.8 | 48.6 | 39.0 | 43.3 | – | 59.3 | 76.6 | 66.8 |
| LBP [73] | 61.9 | 42.7 | 71.6 | 53.5 | 63.2 | 39.2 | 36.9 | 38.1 | – | – | – | – |
| FREE+ESZSL [74] | – | 51.6 | 60.4 | 55.7 | – | 48.2 | 36.5 | 41.5 | – | 51.3 | 78.0 | 61.8 |
| APN+f-VAEGAN-D2 ‡ [75] | 73.9 | 65.5 | 75.6 | <u>70.2</u> | 65.9 | 41.4 | **89.9** | **56.7** | 71.2 | 63.2 | 81.0 | 71.0 |
| *Non-Generative Methods* | | | | | | | | | | | | |
| SP-AEN [76] | 55.4 | 34.7 | 70.6 | 46.6 | 59.2 | 24.9 | 38.6 | 30.3 | 58.5 | 23.3 | 90.9 | 37.1 |
| PQZSL [77] | – | 43.2 | 51.4 | 46.9 | – | 35.1 | 35.3 | 35.2 | – | 31.7 | 70.9 | 43.8 |
| IIR† [78] | 63.8 | 55.8 | 52.3 | 53.0 | 63.5 | 47.9 | 30.4 | 36.8 | 67.9 | 48.5 | 83.2 | 61.3 |
| TCN [79] | 59.5 | 52.6 | 52.0 | 52.3 | 61.5 | 31.2 | 37.3 | 34.0 | 71.2 | 61.2 | 65.8 | 63.4 |
| DVBE [80] | – | 53.2 | 60.2 | 56.5 | – | 45.0 | 37.2 | 40.7 | – | 63.6 | 70.8 | 67.0 |
| DAZLE*† [35] | 66.0 | 56.7 | 59.6 | 58.1 | 59.4 | 52.3 | 24.3 | 33.2 | 67.9 | 60.3 | 75.7 | 67.1 |
| GNDAN*† [81] | 75.1 | <u>69.2</u> | 69.6 | 69.4 | 65.3 | 50.0 | 34.7 | 41.0 | 71.0 | 60.2 | 80.8 | 69.0 |
| MSDN*† [82] | 76.1 | 68.7 | 67.5 | 68.1 | 65.8 | 52.2 | 34.2 | 41.3 | 70.1 | 62.0 | 74.5 | 67.7 |
| TransZero [22](Conference Version) | 76.8 | **69.3** | 68.3 | 68.8 | 65.6 | 52.6 | 33.4 | 40.8 | 70.1 | 61.3 | 82.3 | <u>70.2</u> |
| **TransZero++ (Ours)** | **78.3** | 67.5 | 73.6 | **70.4** | <u>67.6</u> | 48.6 | 37.8 | 42.5 | <u>72.6</u> | 64.6 | 82.7 | **72.5** |

TABLE 2: Conventional zero-shot learning results of various methods on ImageNet. Results indicated with ∗ and † are taken from [83] and [84], respectively.

| Methods | ConSE* [85] | SYNC* [86] | EXEM* [87] | GCNZ* [88] | Trivial† [84] | SGCN* [83] | TransZero++ |
|---|---|---|---|---|---|---|---|
| Top-1 (%) | 8.30 | 10.5 | 12.5 | 19.8 | 20.3 | 26.2 | 23.9 |

ImageNet). To facilitate fair comparison, we follow the train/test split suggested in [83][3]. Table 2 shows the top-1 accuracy of various methods in the "2-hops" setting, which contains all the classes within two hops from the seen classes [15], [83]. We replace attribute features with w2v for TransZero++. Results show that TransZero++ outperforms existing ZSL models. This indicates TransZero++ has potential advantages in ZSL via semantic collaborative learning. Because SGCN [83] employs a graph convolutional neural network (GCN) to discover the prior information of hierarchical relationships among various classes, it achieves better results than our TransZero++.

### 4.1.2 Generalized Zero-Shot Learning

Table 1 shows the results of different methods in the GZSL setting. We can see that most state-of-the-art methods achieve

3. https://github.com/yinboc/DGP

good results on seen classes but fail on unseen classes, while our TransZero++ generalizes better to unseen classes with high unseen and seen accuracies. For example, TransZero++ obtains the best performance with a harmonic mean of 70.4% and 72.5% on CUB and AWA2, respectively. We argue these desirable results benefit from the fact that i) the feature augmentation encoders in AVT and VAT effectively refine the visual features that are more discriminative and transferable than the ones directly extracted from the CNN backbone; ii) the VAT and AVT discover the key semantic knowledge between visual and attribute features for locality-augmented feature representations, enabling effective knowledge transfer from seen to unseen classes. Since Transzero++ is an embedding-based method, it cannot achieve best results on SUN compared to the strong generative methods, *e.g.*, OCD-CVAE [70], TF-VAEGAN [46], HSVA [27] and IZF [29]. Since per class only contains about 16 training images on SUN, which

TABLE 3: Ablation studies for different components of TransZero++ on the CUB and SUN datasets. "FAE" is the feature augmentation encoder, "FA" means feature augmentation, and "DEC" denotes the decoders in AVT and VAT.

| Method | CUB | | | | SUN | | | |
|---|---|---|---|---|---|---|---|---|
| | acc | U | S | H | acc | U | S | H |
| TransZero++ w/o AVT | 49.0 | 36.4 | 48.5 | 41.6 | 57.7 | 36.9 | 28.7 | 32.3 |
| TransZero++ w/o VAT | 75.6 | 62.8 | 72.3 | 67.2 | 63.1 | 44.1 | 35.5 | 39.3 |
| TransZero++ w/o FAE | 70.7 | 63.4 | 57.5 | 60.3 | 64.0 | 52.0 | 33.5 | 40.8 |
| TransZero++ w/o FA | 76.4 | 66.6 | 70.3 | 68.4 | 65.3 | 46.7 | 36.9 | 41.2 |
| TransZero++ w/o DEC | 62.7 | 49.2 | 63.4 | 55.4 | 64.1 | 48.4 | 34.0 | 39.9 |
| TransZero++ (full) | 78.3 | 67.5 | 73.6 | 70.4 | 67.6 | 48.6 | 37.8 | 42.5 |

TABLE 4: Ablation studies for different losses of TransZero++ on the CUB and SUN datasets. Note that $\mathcal{L}_{\text{SCL}} = \mathcal{L}_{\text{f\_SCL}} + \mathcal{L}_{\text{p\_SCL}}$.

| Method | CUB | | | | SUN | | | |
|---|---|---|---|---|---|---|---|---|
| | acc | U | S | H | acc | U | S | H |
| TransZero++(VAT) w/o $\mathcal{L}_{\text{SCL}}$ | 49.0 | 36.4 | 48.5 | 41.6 | 57.7 | 36.9 | 28.7 | 32.3 |
| TransZero++(AVT) w/o $\mathcal{L}_{\text{SCL}}$ | 75.6 | 62.8 | 72.3 | 67.2 | 63.1 | 44.1 | 35.5 | 39.3 |
| TransZero++(VAT) w/ $\mathcal{L}_{\text{SCL}}$ | 49.2 | 37.7 | 51.9 | 43.7 | 63.3 | 48.0 | 31.5 | 38.0 |
| TransZero++(AVT) w/ $\mathcal{L}_{\text{SCL}}$ | 77.6 | 67.2 | 73.4 | 70.2 | 63.8 | 45.3 | 34.7 | 39.3 |
| TransZero++ w/o $\mathcal{L}_{\text{SC}}$ | 77.0 | 46.6 | 76.4 | 58.9 | 65.1 | 41.5 | 36.4 | 38.7 |
| TransZero++ w/o $\mathcal{L}_{\text{AR}}$ | 77.3 | 67.1 | 73.4 | 70.1 | 64.7 | 45.2 | 35.4 | 39.7 |
| TransZero++(AVT and VAT) w/o $\mathcal{L}_{\text{f\_SCL}}$ | 78.1 | 67.9 | 72.1 | 69.9 | 65.6 | 47.4 | 37.5 | 41.9 |
| TransZero++(AVT and VAT) w/o $\mathcal{L}_{\text{d\_SCL}}$ | 75.4 | 64.4 | 71.0 | 67.5 | 65.2 | 46.0 | 37.6 | 41.4 |
| TransZero (full) | 78.3 | 67.5 | 73.6 | 70.4 | 67.6 | 48.6 | 37.8 | 42.5 |

heavily limits the ZSL models, the data augmentation is very effective for improving the performance on SUN. Benefiting from that APN+f-VAEGAN-D2 extracts the discriminate visual features to learn a good transductive generative model (e.g., f-VAEGAN-D2) for data augmentation, it achieves the best performance on SUN. As such, most of the strong generative methods perform better than our TransZero++ (embedding-based method) on SUN. Compared to the latest attention-based method (*e.g.*, APN [39]), our TransZero++ achieves significant improvements of 3.2%, 4.9% and 9.6% in harmonic mean on CUB, SUN and AWA2, respectively. This demonstrates the superiority and great potential of our cross attribute-guided Transformer for the ZSL task. Since the semantic collaborative learning encourages AVT and VAT to learn semantic-augmented embedding for desirable visual-semantic interaction, TransZero++ continuously improves the performance of its conference version (TransZero [22]) on all datasets.

## 4.2 Ablation Study

To provide further insight into TransZero++, we conduct ablation studies to evaluate the effect of different model components, loss functions, and distance metrics for semantical collaborative losses. **Analysis of Model Components.** As shown in Table 3, we conduct ablation studies to evaluate the effects of different model components, *i.e.*, AVT, VAT, feature augmentation encoder (denoted as FAE), feature augmentation in FAE (denoted as FA), and visual-semantic decoder (denoted as DEC). We observed that TransZero++ with various model components achieve more clear improvements on CUB than on SUN, because SUN is a scene dataset, which is more complicated and challenging than CUB. TransZero++ performs significantly worse than if no AVT is used, *i.e.*, the acc/harmonic mean drops by 29.3%/28.8% on CUB and 9.9%/10.2% on SUN. This indicates that AVT is an basic attention sub-net in TransZero++, which is consistent with

most existing attention-based methods [35], [39], [51], [91] based on attribute→visual attention. Meanwhile, Transzero++ w/o VAT also achieves inferior performance than the full model. As such, it is necessary to simultaneously learn the semantic-augmented embeddings with AVT and VAT for ZSL, under the guidance of semantical collaborative learning. TransZero++ significantly improves its performance when AVT and VAT use the feature augmentation encoders, which shows the importance of refining the visual feature to alleviate the cross-dataset bias. If we incorporate the encoder of the standard Transformer without feature augmentation, TransZero++ obtains inferior performances as the entangled relative geometry priors limit the transferability of visual features, *i.e.*, the acc/harmonic mean drops by 1.9%/2.0% and 2.3%/1.3% on CUB and SUN, respectively. When TransZero++ does not employ the decoders in AVT and VAT to localize the key object attribute for semantic knowledge representations, its performance decreases dramatically on all datasets.

**Analysis of Loss Functions.** As shown in Table 4, we further conduct ablation studies to evaluate the effects of different loss functions, *i.e.*, semantic collaborative loss (including $\mathcal{L}_{\text{f\_SCL}}$ and $\mathcal{L}_{\text{p\_SCL}}$), self-calibration loss (*i.e.*, $\mathcal{L}_{\text{SC}}$) and attribute regression loss (*i.e.*, $\mathcal{L}_{\text{AR}}$). TransZero++ using single sub-net (*i.e.*, TransZero++(VAT) and TransZero++(VAT)) achieves significant gains on CUB and SUN when it uses the semantical collaborative loss ($\mathcal{L}_{\text{SCL}} = \mathcal{L}_{\text{f\_SCL}} + \mathcal{L}_{\text{p\_SCL}}$). For example, TransZero++(VAT) and TransZero++(AVT) achieve gains of 2.1% and 3.0% in harmonic mean on CUB, respectively. This shows that semantic collaborative learning is effective for encouraging AVT and VAT to teach each other to discover the key transferable semantic knowledge for ZSL. The self-calibration mechanism can effectively alleviate the seen-unseen bias problem [35], [38], [39], resulting in improvements in the harmonic mean of 11.5% on CUB. The attribute regression constraint further improves the performance of TransZero++ by directing $\mathcal{M}_1$ and $\mathcal{M}_2$ to conduct effective visual-semantic

TABLE 5: Ablation studies for different components of TransZero++ on the CUB and SUN datasets. "FAE" is the feature augmentation encoder, "FA" means feature augmentation, and "DEC" denotes visual-semantic decoder.

| Method | CUB | | | | SUN | | | |
|---|---|---|---|---|---|---|---|---|
| | acc | U | S | H | acc | U | S | H |
| TransZero++ w/ $\ell_1$ | 53.6 | 42.6 | 48.9 | 37.8 | 64.3 | 51.4 | 32.7 | 40.0 |
| TransZero++ w/ KL($p_1\|\|p_2$) | 76.4 | 67.8 | 65.4 | 70.4 | 65.0 | 46.3 | 37.6 | 41.5 |
| TransZero++ w/ KL($p_2\|\|p_1$) | 76.3 | 68.3 | 65.3 | 71.6 | 65.0 | 46.3 | 38.1 | 41.8 |
| TransZero++ w/ JSD | 76.1 | 67.2 | 63.2 | 71.8 | 64.8 | 46.2 | 37.9 | 41.6 |
| TransZero++ w/ $\ell_2$ | 78.3 | 67.5 | 73.6 | 70.4 | 67.6 | 48.6 | 37.8 | 42.5 |

mapping. Furthemore, the two semantic collaborative losses (*i.e.*, $\mathcal{L}_{\text{f\_SCL}}$ and $\mathcal{L}_{\text{p\_SCL}}$) encourage TransZero++ to conduct desirable semantic collaborative learning.

**Analysis of Distance Metrics for Semantical Collaborative Losses.** As shown in Table 5, we conduct ablation studies to evaluate the effects of distance metrics for semantic collaborative losses (feature-level semantic collaborative loss ($\mathcal{L}_{\text{f\_SCL}}$) and prediction-level semantic collaborative loss ($\mathcal{L}_{\text{p\_SCL}}$)), *i.e.*, $\ell_1$, $\ell_2$, KL($p_1\|\|p_2$), KL($p_2\|\|p_1$), and JSD. Results show that TransZero++ performs very poorly using $\ell_1$ distance for calculating the semantic collaborative losses. The possible reason is that the values of semantic-augmented visual embeddings and predictions in VAT and AVT are constrained to be in a small range, $\ell_1$ distance cannot well capture the divergence between the outputs of AVT and VAT. When TransZero++ use the KL($p_1\|\|p_2$), KL($p_2\|\|p_1$), or JSD to compute the semantic collaborative losses, it achieves consistent good performance almost. Thus, the symmetric and asymmetric distances for semantic collaborative losses do not make any difference. Interestingly, TransZero++ achieves the best results using $\ell_2$ distance which is beneficial for the regression problem. As such, we take $\ell_2$ to compute $\mathcal{L}_{\text{p\_SCL}}$ and $\mathcal{L}_{\text{f\_SCL}}$.

## 4.3 Qualitative Results

Here, we present the visualizations of attention maps and t-SNE [92] to intuitively show the effectiveness of our TransZero++.

### 4.3.1 Visualization of Attention Maps.

To intuitively show the effectiveness of our TransZero++ at learning attribute-relevant visual features, we visualize the attention maps learned by the existing attention-based methods (*e.g.*, AREN [37]) and TransZero++. As shown in Fig. 5, AREN simply learns region embeddings for visual representations, *e.g.*, the whole bird body, neglecting the fine-grained semantic attribute information. In contrast, our Transzero++ learns discriminative attribute localization for visual features by assigning high positive scores to key attributes (*e.g.*, the "bill shape all-purpose" of the *Acadian Flycatcher* in Fig. 5). Thus, TransZero++ discovers the semantic-augmented embeddings that are discriminative and transferable, enabling good performance both in seen and unseen classes. Compared to TransZero [22] (Conference version), TransZero++ can discover more valuable attributes for semantic-augmented embedding representations (*e.g.*, "upper part color gray" of *Acadian Flycatcher*). Furthermore, TransZero++ gets higher confidence scores for the important attributes that exist in the images than TransZero. For example, TransZero++ gets the scores of 28.0 for attribute "head pattern plain" of the *Acadian Flycatcher*, while TransZero gets the scores of 14.7.

### 4.3.2 t-SNE Visualizations.

As shown in Fig. 6, we provide the t-SNE visualization [92] of visual features for (a) seen classes and (b) unseen classes on CUB, learned by the CNN backbone, TransZero++(AVT) encoder w/o FA, TransZero++(AVT) encoder, TransZero++(AVT) decoder, TransZero++(VAT) encoder w/o FA, TransZero++(VAT) encoder, TransZero++(VAT) decoder, TransZero++(AVT and VAT). If the standard encoder is incorporated into AVT and VAT of our TransZero++, the visual features learned by the encoder are significantly improved compared to the original visual features extracted from the CNN Backbone (*e.g.*, ResNet101). When we use the feature augmentation encoder to refine the original visual features, the quality of visual features is further enhanced. These results demonstrate that the encoder of TransZero++ effectively alleviates the cross-dataset bias problem and reduces the entangled relative geometry relationships among different regions, enabling the visual feature to be more discriminative and transferable. Moreover, the attribute→visual and visual→attribute decoders in AVT and VAT learn attribute-based visual features and visual-based attribute features, which are further mapped into semantic embedding space for semantic-augmented embedding representations. Since the features learned by AVT and VAT are complementary to each other, the fused semantic-augmented embedding can be further refined. As such, our TransZero++ achieves significant performance both in seen and unseen classes on all datasets.

## 4.4 Part Localization Prediction

Together with attention map visualizations, we quantitatively report the part localization prediction by calculating the Percentage of Correctly Localized Parts (PCP) of unseen classes on CUB. Our calculation follows that in APN[4] [39]. When the predicted bounding box for a part overlaps sufficiently with the grounding truth bounding box (i.e., IoU > 0.5), the detection is considered to be correct. Table 6 shows the results where our TransZero++ consistently and significantly improves the localization accuracy of all parts over DAZLE [35] and APN [39]. These results are in accordance to the qualitative results in Fig. 5 where our TransZero++ can better discover the key semantic knowledge between visual and attribute features, resulting in desirable knowledge transfer for ZSL.

TABLE 6: Results of part localization prediction for various methods on CUB. For BB (bounding box) size, $1/\sqrt{2}$ denotes each part bounding box has the size of $1/\sqrt{2}W \times 1/\sqrt{2}H$, where $W$ and $H$ are the width and height of the bird image.

| Methods | BB size | Head | Breast | Belly | Back | Wing | Leg | Mean |
|---|---|---|---|---|---|---|---|---|
| DAZLE [35] | $1/\sqrt{2}$ | 94.7 | 92.1 | 75.1 | 65.1 | 68.4 | 50.5 | 74.3 |
| APN [39] | $1/\sqrt{2}$ | 91.8 | 88.9 | 81.0 | 72.1 | 76.6 | 65.0 | 79.2 |
| TransZero++ | $1/\sqrt{2}$ | 97.9 | 93.7 | 85.3 | 80.4 | 86.0 | 74.0 | 86.2 |

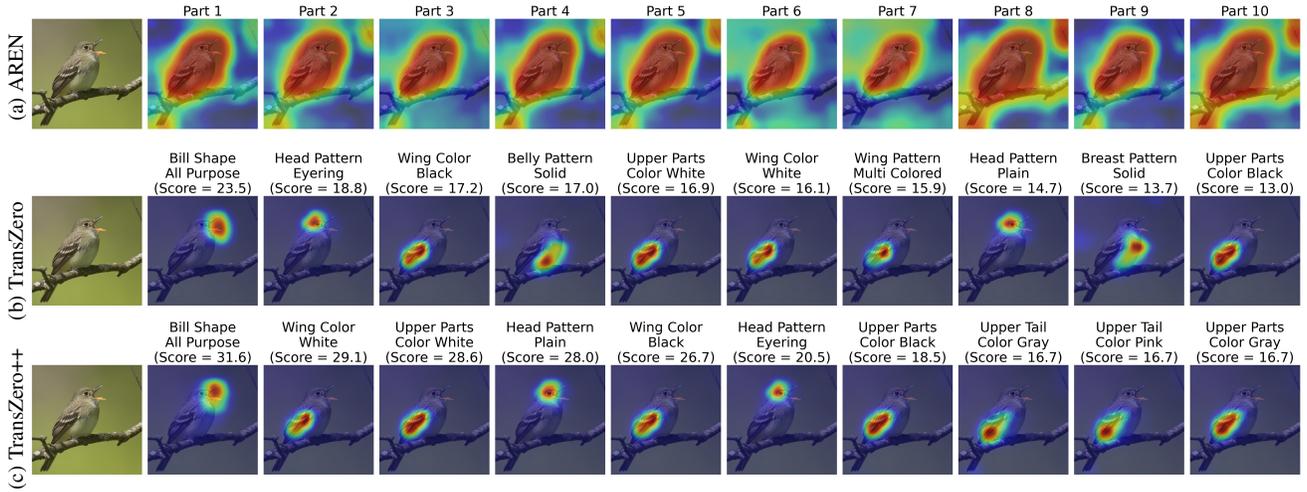4. https://github.com/wenjiaXu/APN-ZSL

Fig. 5: Visualization of top-10 attention maps for the attention-based method (*i.e.*, AREN [37]), TransZero [22] (Conference version) and our TransZero++. Results show that TransZero localizes some important object attributes with low confidence scores for representing region features, while AREN is failed. Furthermore, our TransZero++ discovers more valuable attributes that exist in the corresponding image with high confidence scores compared to the TransZero. More results are presented in the Project Website. (Best viewed in color)

## 4.5 Generality Analysis

To show the generality of our TransZero++, we conduct experiments by replacing the attributes descriptions with word vectors of classes name (denotes w2v) [93]. Results are shown in Table 7. Since the attribute descriptions provide more informative representations than w2v, all ZSL methods achieve inferior performance using w2v compared to attributes. Fortunately, we observe that our method achieves the best result of 25.1% and 28.7% in the CZSL setting on CUB and SUN, respectively. In the GZSL setting, TransZero++ performs significant gains of harmonic mean with 5.1% over APN [39] on SUN dataset. These results indicate that TransZero++ can also discover the key semantic knowledge with the guidance of w2v. This further demonstrates the advantages of our method.

TABLE 7: Results of various methods with word vector of class names (denotes as w2v). "†" indicates that results are taken from [94].

| Method | CUB | | | | SUN | | | |
|---|---|---|---|---|---|---|---|---|
| | acc | U | S | H | acc | U | S | H |
| SJE† [95] | 14.4 | 13.2 | 28.6 | 18.0 | 26.3 | 19.8 | 18.6 | 19.2 |
| APN† [39] | 22.7 | 17.6 | 29.4 | 22.1 | 23.6 | 16.3 | 15.3 | 15.8 |
| TransZero++ | 25.1 | 12.9 | 29.9 | 18.0 | 28.7 | 16.4 | 28.7 | 20.9 |

Furthermore, we also extract feature representations from the decoders of our TransZero++ and entailing them on the top of generative models, *e.g.*, f-VAEGAN [8] and TF-VAEGAN [46]. Results in Table 8 show that our TransZero++ consistently boosts the performances of generative models on two datasets. Specifically, TransZero++ improves the performance of f-VAEGAN on $acc/H$ with 10.2%/7.1% and 1.4%/0.3% on CUB and SUN, respectively. TransZero++ also gains the improvements of TF-VAEGAN on $acc/H$ with 6.6%/5.0% and 2.5%/1.3% on CUB and SUN, respectively. These results indicate that our TransZero++ discovers the key semantic knowledge to represent the semantic-augmented features, helping the generative models synthesize discriminative and transferable visual features. As such, the cross-dataset bias problem [28] in f-VAEGAN and TF-VAEGAN can be alleviated.

Finally, we also extend TransZero++ to the transductive ZSL setting following [18], where unlabeled samples of unseen classes are also used for model optimization. We compare TransZero++

TABLE 8: Results of various generative models with visual features extracted from TransZero++.

| Method | CUB | | | | SUN | | | |
|---|---|---|---|---|---|---|---|---|
| | acc | U | S | H | acc | U | S | H |
| f-VAEGAN [8] | 61.0 | 48.4 | 60.1 | 53.6 | 64.7 | 45.1 | 38.0 | 41.3 |
| TransZero++ ⊕ f-VAEGAN [8] | 71.2 | 61.7 | 59.7 | 60.7 | 66.1 | 50.4 | 35.5 | 41.6 |
| TF-VAEGAN [46] | 64.9 | 52.8 | 64.7 | 58.1 | 66.0 | 45.6 | 40.7 | 43.0 |
| TransZero++ ⊕ TF-VAEGAN [46] | 71.5 | 63.3 | 62.8 | 63.1 | 68.5 | 48.3 | 40.9 | 44.3 |

with other embedding-based Transductive ZSL methods on CUB and SUN datasets, as shown in Table 9. Compared to other embedding-based transductive ZSL methods, our TransZero++ achieves new state-of-the-art of $acc/H$ with 81.5%/77.7% and 69.0%/48.8% on CUB and SUN, respectively. This indicates that TransZero++ also learns the intrinsic semantic knowledge for desirable knowledge transfer in transductive ZSL.

TABLE 9: Results of various embedding-based transductive ZSL methods on CUB and SUN datasets.

| Method | CUB | | | | SUN | | | |
|---|---|---|---|---|---|---|---|---|
| | acc | U | S | H | acc | U | S | H |
| ALE-tran [15] | 54.5 | 23.5 | 45.1 | 30.9 | 55.7 | 19.9 | 22.6 | 21.2 |
| DSRL [96] | 48.7 | 17.3 | 39.0 | 24.0 | 56.8 | 17.7 | 25.0 | 20.7 |
| UE-finetune [18] | 72.1 | 74.9 | 71.5 | 73.2 | 58.3 | 33.6 | 54.8 | 41.7 |
| TransZero++ | 81.5 | 79.6 | 75.8 | 77.7 | 69.0 | 64.9 | 39.1 | 48.8 |

## 4.6 Hyperparameter Analysis

To analyse the robustness of our TransZero++ and select better hyperparameters for it. We conduct extensive experiments for evaluating the effects of loss weights (in Eq. 29 and Eq. 31), Transformer architecture settings in AVT and VAT, and combination coefficient (in Eq. 32).

### 4.6.1 Effects of Loss Weights

Here, we analyse the effects of loss weights that control their corresponding loss terms, *i.e.*, $\lambda_{\text{AR}}$, $\lambda_{\text{SC}}$, $\lambda_{\text{VAT}}$, $\lambda_{\text{f\_SCL}}$ and $\lambda_{\text{p\_SCL}}$. We try a range of these loss weights evaluated on CUB and SUN, *i.e.*, $\{0.0001, 0.001, 0.01, 0.1, 1.0\}$. Results are shown in Fig. 7. When $\lambda_{AR}$, $\lambda_{\text{f\_SCL}}$ and $\lambda_{\text{p\_SCL}}$ are set to a large value, all evaluation protocols tend to drop. Moreover, TansZero++ are relatively insensitive to $\lambda_{\text{SC}}$ and $\lambda_{VAT}$ when they are set to small (*e.g.*, smaller than 0.01). Based on
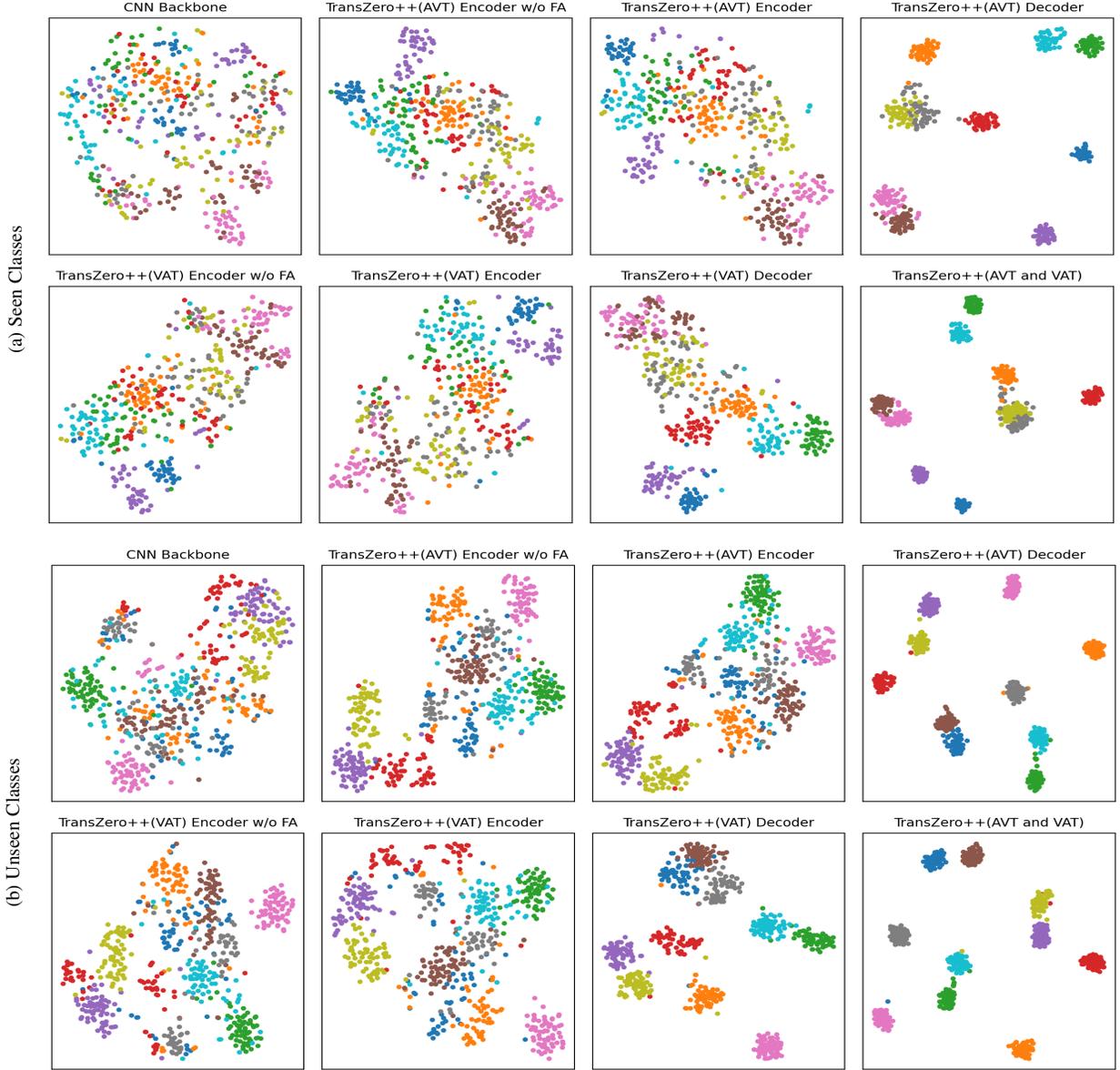
Fig. 6: t-SNE visualizations of visual features for (a) seen classes and (b) unseen classes, learned by the CNN backbone, TransZero++(AVT) encoder w/o FA, TransZero++(AVT) encoder, TransZero++(AVT) decoder, TransZero++(VAT) encoder w/o FA, TransZero++(VAT) encoder, TransZero++(VAT) decoder and TransZero++(VAT and VAT). The 10 colors denote 10 different seen/unseen classes randomly selected from CUB. Results show that our various model components in TransZero++ learn the discriminative visual feature representations, while CNN backbone (*e.g.*, ResNet101) failed. The results on SUN and AWA2 are presented in the Project Website. (Best viewed in color)

these observations, we set $\{\lambda_{\mathrm{AR}}, \lambda_{\mathrm{SC}}, \lambda_{\mathrm{VAT}}, \lambda_{\mathrm{f\_SCL}}, \lambda_{\mathrm{p\_SCL}}\}$ to $\{0.0001, 0.1, 0.1, 0.001, 0.01\}$ and $\{0.01, 0.1, 1.0, 0.001, 0.001\}$ for CUB and SUN datasets, repectively.

### 4.6.2 Effects of Different Architectures for Transformer in AVT and VAT

To find the best Transformer settings in AVT and VAT, we investigate the influence of the number of i) layers of Encoder/Decoder, and ii) attention heads of Encoder/Decoder. To enable the training of TrasZero++ to be more stable, we set same number of layers/heads in the encoder and decoder. As shown in Fig. 8, we find that the encoders/decoders in AVT and VAT should be set to be small, *i.e.*, one layer with one attention head, TransZero++ can achieve better results in both seen and unseen classes. The

possible reason lies in that the training data for the ZSL model is medium/small scale which inevitably leads to over-fitting with more complex Transformer architectures.

### 4.6.3 Effects of Combination Coefficient

We argue that the attribute-based visual features and visual-based attribute features learned by AVT and VAT respectively are complementary, and thus we take a combination coefficient $\alpha$ to fuse their corresponding semantic-augmented embeddings for desirable visual-semantic interaction (in Eq. 32). We try a range of $\alpha$ on CUB and SUN, *i.e.*, $\{0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$. Notably, $\alpha = 0.0$ is denoted as the TransZero++(VAT), and $\alpha = 1.0$ is denoted as the TransZero++(AVT). As shown in Fig. 9, when $\alpha$ is set to large relatively (*e.g.*, $\alpha > 0.5$), TransZero++
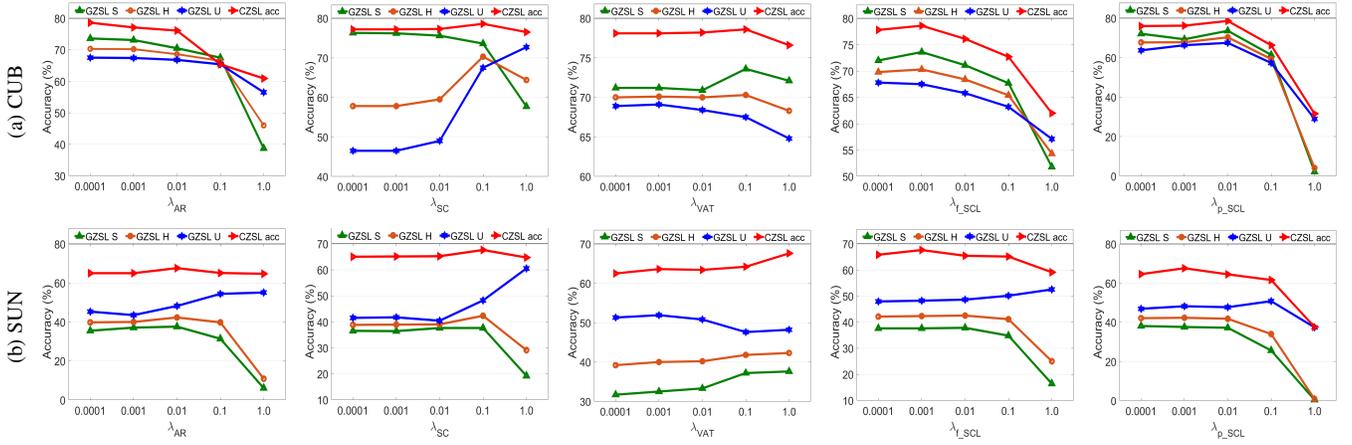
Fig. 7: The effects of loss weights that control their corresponding loss terms on CUB and SUN, *i.e.*, $\lambda_{AR}$, $\lambda_{SC}$, $\lambda_{VAT}$, $\lambda_{f\_SCL}$ and $\lambda_{p\_SCL}$.
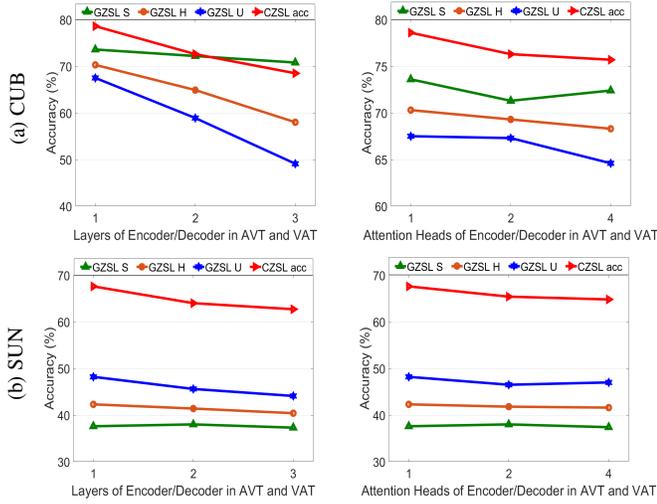


Fig. 8: The effects of different architectures for the AVT and VAT networks on CUB and SUN. We investigate the number of layers of Encoder/Decoder and attention heads in Encoder/Decoder.

achieves better results. This demonstrates that AVT sub-net provides more desirable information for TransZero++. However, $\alpha$ should also not be set too large since the VAT sun-net provides additional useful information for TransZero++. Based on these results, we set $\alpha$ to 0.9 and 0.6 for CUB and SUN, respectively.

## 5 CONCLUSION

In this paper, we propose a novel cross attribute-guided Transformer network for ZSL, termed TransZero++. TransZero++ consists of a attribute→visual Transformer sub-net (AVT) and visual→attribute Transformer sub-net (VAT). First, AVT employs a feature augmentation encoder to improve the discriminability and transferability of visual features by alleviating the cross-dataset problem and reducing the entangled region feature relationships, respectively. Meanwhile, an attribute→visual decoder in AVT is introduced to learn the attribute localization for attribute-based visual feature representations which are locality-augmented. Secondly, VAT applied a similar feature augmentation encoder to refine the visual features, which is further fed into a visual→attribute decoder to learn the visual-based attribute features. By introducing the feature-level and prediction-level semantical collaborative losses for optimization, our TransZero++ can learn the semantic-augmented
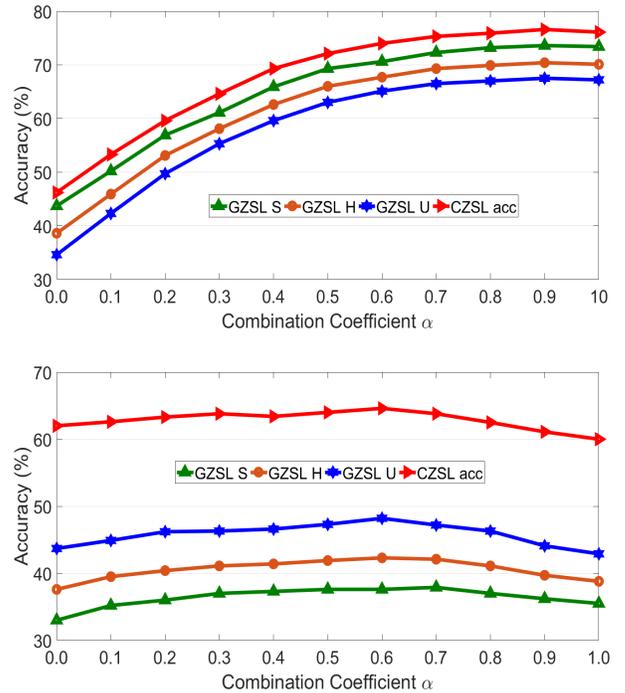


Fig. 9: The effectiveness of combination coefficients $\alpha$ between the AVT and VAT sub-nets on CUB (top) and SUN (bottom).

visual embedding. Considering the attribute-based visual features and visual-based attribute features that are complementary to each other, we combine the two semantic-augmented visual embeddings learned by AVT and VAT to enable desirable visual-semantic interaction cooperated with the class semantic vectors for ZSL classification. Extensive experiments on three popular ZSL benchmarks and on the large-scale ImageNet dataset demonstrate the superiority of our method. We believe that our work also facilitates the development of other visual-and-language learning systems, *e.g.*, image captioning [69], natural language for visual reasoning [97].

Indeed, our TransZero++ bases on the Glove embedding of attribute names, which is not easy to get in real-life. Since w2v lacks enough informative representations, it cannot well supports ZSL methods to significantly discover the key semantic knowledge [94], [98]. Motivated by the work of Xu [94], we can effectively discover semantic embeddings containing discriminative visual

properties via visually clustering and class relations prediction, without requiring any human annotation to replace the attributes defined by experts in future work.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] H. Larochelle, D. Erhan, and Y. Bengio, "Zero-data learning of new tasks," in *AAAI*, 2008, pp. 646–651.

[2] M. Palatucci, D. Pomerleau, G. E. Hinton, and T. M. Mitchell, "Zero-shot learning with semantic output codes," in *NeurIPS*, 2009, pp. 1410–1418.

[3] C. H. Lampert, H. Nickisch, and S. Harmeling, "Learning to detect unseen object classes by between-class attribute transfer," in *CVPR*, 2009, pp. 951–958.

[4] C. H. Lampert, S. Harmeling, and H. Nickisch, "Attribute-based classification for zero-shot visual object categorization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, pp. 453–465, 2014.

[5] Z. Fu, T. Xiang, E. Kodirov, and S. Gong, "Zero-shot learning on semantic class prototype graph," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, pp. 2009–2022, 2018.

[6] FuYanwei, M. HospedalesTimothy, Xiang-tao, and GongShaogang, "Transductive multi-view zero-shot learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 2332–2345, 2015.

[7] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A deep visual-semantic embedding model," in *NeurIPS*, 2013.

[8] Y. Xian, S. Sharma, B. Schiele, and Z. Akata, "F-vaegan-d2: A feature generating framework for any-shot learning," in *CVPR*, 2019, pp. 10 267–10 276.

[9] Y. Shen, L. Liu, F. Shen, and L. Shao, "Zero-shot sketch-image hashing," in *CVPR*, 2018, pp. 3598–3607.

[10] A. Dutta and Z. Akata, "Semantically tied paired cycle consistency for any-shot sketch-based image retrieval," *International Journal of Computer Vision*, pp. 1–20, 2020.

[11] M. Bucher, T.-H. Vu, M. Cord, and P. Pérez, "Zero-shot semantic segmentation," in *NeurIPS*, 2019, pp. 468–479.

[12] A. Bansal, K. Sikka, G. Sharma, R. Chellappa, and A. Divakaran, "Zero-shot object detection," in *ECCV*, 2018.

[13] S. Reed, Z. Akata, H. Lee, and B. Schiele, "Learning deep representations of fine-grained visual descriptions," in *CVPR*, 2016, pp. 49–58.

[14] S. Badirli, Z. Akata, G. O. Mohler, C. Picard, and M. Dundar, "Fine-grained zero-shot learning with dna as side information," in *NeurIPS*, 2021.

[15] Y. Xian, C. H. Lampert, B. Schiele, and Z. Akata, "Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, pp. 2251–2265, 2019.

[16] Y. Xian, T. Lorenz, B. Schiele, and Z. Akata, "Feature generating networks for zero-shot learning," in *CVPR*, 2018, pp. 5542–5551.

[17] G.-S. Xie, Z. Zhang, G.-S. Liu, F. Zhu, L. Liu, L. Shao, and X. Li, "Generalized zero-shot learning with multiple graph adaptive generative networks." *IEEE transactions on neural networks and learning systems*, 2021.

[18] J. Song, C. Shen, Y. Yang, Y. Liu, and M. Song, "Transductive unbiased embedding for zero-shot learning," *CVPR*, pp. 1024–1033, 2018.

[19] G.-S. Xie, X.-Y. Zhang, Y. Yao, Z. Zhang, F. Zhao, and L. Shao, "Vman: A virtual mainstay alignment network for transductive zero-shot learning," *IEEE Transactions on Image Processing*, vol. 30, pp. 4316–4329, 2021.

[20] W. Yin, J. Hay, and D. Roth, "Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach," in *EMNLP*, 2019.

[21] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.

[22] S. Chen, Z. Hong, Y. Liu, G.-S. Xie, B. Sun, H. Li, Q. Peng, K. Lu, and X. You, "Transzero: Attribute-guided transformer for zero-shot learning," in *AAAI*, 2022.

[23] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid, "Label-embedding for image classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, pp. 1425–1438, 2016.

[24] Y. Li, J. Zhang, J. Zhang, and K. Huang, "Discriminative learning of latent features for zero-shot recognition," in *CVPR*, 2018, pp. 7463–7471.

[25] G. Arora, V. Verma, A. Mishra, and P. Rai, "Generalized zero-shot learning via synthesized examples," in *CVPR*, 2018, pp. 4281–4289.

[26] E. Schönfeld, S. Ebrahimi, S. Sinha, T. Darrell, and Z. Akata, "Generalized zero- and few-shot learning via aligned variational autoencoders," in *CVPR*, 2019, pp. 8239–8247.

[27] S. Chen, G.-S. Xie, Y. Yang Liu, Q. Peng, B. Sun, H. Li, X. You, and L. Shao, "Hsva: Hierarchical semantic-visual adaptation for zero-shot learning," in *NeurIPS*, 2021.

[28] S. Chen, W. Wang, B. Xia, Q. Peng, X. You, F. Zheng, and L. Shao, "Free: Feature refinement for generalized zero-shot learning," in *ICCV*, 2021.

[29] Y. Shen, J. Qin, and L. Huang, "Invertible zero-shot recognition flows," in *ECCV*, 2020.

[30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[31] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255.

[32] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. J. Belongie, and P. Perona, "Caltech-ucsd birds 200," *Technical Report CNS-TR-2010-001, Caltech,*, 2010.

[33] Z. Chen, Y. Luo, R. Qiu, S. Wang, Z.-Y. Huang, J. Li, and Z. Zhang, "Semantics disentangling for generalized zero-shot learning," in *ICCV*, 2021.

[34] G.-S. Xie, L. Liu, X. Jin, F. Zhu, Z. Zhang, Y. Yao, J. Qin, and L. Shao, "Region graph embedding network for zero-shot learning," in *ECCV*, 2020.

[35] D. Huynh and E. Elhamifar, "Fine-grained generalized zero-shot learning via dense attribute-based attention," in *CVPR*, 2020, pp. 4482–4492.

[36] Z. Wang, Y. Gou, J. Li, Y. Zhang, and Y. Yang, "Region semantically aligned network for zero-shot learning," in *CIKM*, 2021.

[37] G.-S. Xie, L. Liu, X. Jin, F. Zhu, Z. Zhang, J. Qin, Y. Yao, and L. Shao, "Attentive region embedding network for zero-shot learning," in *CVPR*, 2019, pp. 9376–9385.

[38] Y. Zhu, J. Xie, Z. Tang, X. Peng, and A. Elgammal, "Semantic-guided multi-attention localization for zero-shot learning," in *NeurIPS*, 2019.

[39] W. Xu, Y. Xian, J. Wang, B. Schiele, and Z. Akata, "Attribute prototype network for zero-shot learning," in *NeurIPS*, 2020.

[40] Y. Yu, Z. Ji, Y. Fu, J. Guo, Y. Pang, and Z. Zhang, "Stacked semantics-guided attention model for fine-grained zero-shot learning," in *NeurIPS*, 2018.

[41] Y. Liu, J. Guo, D. Cai, and X. He, "Attribute attention for semantic disambiguation in zero-shot learning," in *ICCV*, 2019, pp. 6697–6706.

[42] Z. Hong, S. Chen, G. Xie, W. Yang, J. Zhao, Y. Shao, Q. Peng, and X. You, "Semantic compression embedding for generative zero-shot learning," in *IJCAI*, 2022.

[43] G. Patterson and J. Hays, "Sun attribute database: Discovering, annotating, and recognizing scene attributes," in *CVPR*, 2012, pp. 2751–2758.

[44] Z. Han, Z. Fu, S. Chen, and J. Yang, "Contrastive embedding for generalized zero-shot learning," in *CVPR*, 2021.

[45] Y.-Y. Chou, H.-T. Lin, and T.-L. Liu, "Adaptive and generative zero-shot learning," in *ICLR*, 2021.

[46] S. Narayan, A. Gupta, F. Khan, C. G. M. Snoek, and L. Shao, "Latent embedding feedback and discriminative features for zero-shot classification," in *ECCV*, 2020.

[47] C. Yan, X. Chang, Z. Li, Z. Ge, W. Guan, L. Zhu, and Q. Zheng, "Zeronas: Differentiable generative adversarial networks search for zero-shot learning." *IEEE transactions on pattern analysis and machine intelligence*, 2021.

[48] Y.-H. H. Tsai, L.-K. Huang, and R. Salakhutdinov, "Learning robust visual-semantic embeddings," in *ICCV*, 2017, pp. 3591–3600.

[49] Y. Atzmon, F. Kreuk, U. Shalit, and G. Chechik, "A causal view of compositional zero-shot recognition," in *NeurIPS*, 2020.

[50] T. Chen, T. Pu, Y. Xie, H. Wu, L. Liu, and L. Lin, "Cross-domain facial expression recognition: A unified evaluation benchmark and adversarial graph learning." *IEEE transactions on pattern analysis and machine intelligence*, 2021.

[51] Y. Liu, L. Zhou, X. Bai, Y. Huang, L. Gu, J. Zhou, and T. Harada, "Goal-oriented gaze estimation for zero-shot learning," in *CVPR*, 2021.

[52] A. Vaswani, N. M. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017.

[53] M. Naseer, K. Ranasinghe, S. H. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Intriguing properties of vision transformers," in *NeurIPS*, 2021.

[54] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. Khan, and M. Shah, "Transformers in vision: A survey," *arXiv preprint arXiv:2101.01169*, 2021.

[55] W. Yu, M. Luo, P. Zhou, C. Si, Y. Zhou, X. Wang, J. Feng, and S. Yan, "Metaformer is actually what you need for vision," *arXiv preprint arXiv:2111.11418*, 2021.

[56] M. Ott, S. Edunov, D. Grangier, and M. Auli, "Scaling neural machine translation," in *WMT*, 2018.

[57] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021.

[58] V. Gabeur, C. Sun, A. Karteek, and C. Schmid, "Multi-modal transformer for video retrieval," in *ECCV*, 2020.

[59] X. Zhang, X. Sun, Y. Luo, J. Ji, Y. Zhou, Y. Wu, F. Huang, and R. Ji, "Rstnet: Captioning with adaptive attention on visual and non-visual words," in *CVPR*, 2021.

[60] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara, "Meshed-memory transformer for image captioning," in *CVPR*, 2020, pp. 10 575–10 584.

[61] L. Huang, W. Wang, J. Chen, and X.-Y. Wei, "Attention on attention for image captioning," in *ICCV*, 2019, pp. 4633–4642.

[62] Y. Pan, T. Yao, Y. Li, and T. Mei, "X-linear attention networks for image captioning," in *CVPR*, 2020, pp. 10 968–10 977.

[63] T. Batra and D. Parikh, "Cooperative learning with visual attributes," *arXiv preprint arXiv: 1705.05512*, 2017.

[64] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *NeurIPS*, 2017.

[65] Y. Ge, D. peng Chen, and H. Li, "Mutual mean-teaching: Pseudo label refinery for unsupervised domain adaptation on person re-identification," in *ICLR*, 2020.

[66] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," *CVPR*, pp. 4320–4328, 2018.

[67] Y. Zhai, Q. Ye, S. Lu, M. Jia, R. Ji, and Y. Tian, "Multiple expert brainstorming for domain adaptive person re-identification," in *ECCV*, 2020.

[68] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014.

[69] S. Herdade, A. Kappeler, K. Boakye, and J. Soares, "Image captioning: Transforming objects into words," in *NeurIPS*, 2019.

[70] R. Keshari, R. Singh, and M. Vatsa, "Generalized zero-shot learning via over-complete distribution," in *CVPR*, 2020, pp. 13 297–13 305.

[71] D. T. Huynh and E. Elhamifar, "Compositional zero-shot learning via fine-grained dense feature composition," in *NeurIPS*, 2020.

[72] Z. Yue, T. Wang, H. Zhang, Q. Sun, and X. Hua, "Counterfactual zero-shot and open-set visual recognition," in *CVPR*, 2021.

[73] Z. Lu, J. Guan, A. Li, T. Xiang, A. Zhao, and J.-R. Wen, "Zero and few shot learning with semantic feature synthesis and competitive learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, pp. 2510–2523, 2021.

[74] S. Cetin, O. B. Baran, and R. G. Cinbis, "Closed-form sample probing for learning generative models in zero-shot learning," in *ICLR*, 2022.

[75] W. Xu, Y. Xian, J. Wang, B. Schiele, and Z. Akata, "Attribute prototype network for any-shot learning," *International Journal of Computer Vision*, vol. 130, p. 1735–1753, 2022.

[76] L. Chen, H. Zhang, J. Xiao, W. Liu, and S. Chang, "Zero-shot visual recognition using semantics-preserving adversarial embedding networks," in *CVPR*, 2018, pp. 1043–1052.

[77] J. Li, X. Lan, Y. Liu, L. Wang, and N. Zheng, "Compressing unknown images with product quantizer for efficient zero-shot classification," in *CVPR*, 2019, pp. 5458–5467.

[78] Y. L. Cacheux, H. Borgne, and M. Crucianu, "Modeling inter and intra-class relations in the triplet loss for zero-shot learning," in *ICCV*, 2019, pp. 10 332–10 341.

[79] H. Jiang, R. Wang, S. Shan, and X. Chen, "Transferable contrastive network for generalized zero-shot learning," in *ICCV*, 2019, pp. 9764–9773.

[80] S. Min, H. Yao, H. Xie, C. Wang, Z. Zha, and Y. Zhang, "Domain-aware visual bias eliminating for generalized zero-shot learning," in *CVPR*, 2020, pp. 12 661–12 670.

[81] S. Chen, Z. Hong, G. Xie, Q. Peng, X. You, W. Ding, and L. Shao, "Gndan: Graph navigated dual attention network for zero-shot learning." *IEEE transactions on neural networks and learning systems*, 2022.

[82] S. Chen, Z. Hong, G. Xie, W. Wang, Q. Peng, K. Wang, J. Zhao, and X. You, "Msdn: Mutually semantic distillation network for zero-shot learning," in *CVPR*, 2022, pp. 7612–7621.

[83] M. C. Kampffmeyer, Y. Chen, X. Liang, H. Wang, Y. Zhang, and E. Xing, "Rethinking knowledge graph propagation for zero-shot learning," in *CVPR*, 2019, pp. 11 479–11 488.

[84] T. Hascoet, Y. Ariki, and T. Takiguchi, "On zero-shot recognition of generic objects," in *CVPR*, 2019, pp. 9545–9553.

[85] M. Norouzi, T. Mikolov, S. Bengio, Y. Singer, J. Shlens, A. Frome, G. S. Corrado, and J. Dean, "Zero-shot learning by convex combination of semantic embeddings," in *ICLR*, 2014.

[86] S. Changpinyo, W.-L. Chao, B. Gong, and F. Sha, "Synthesized classifiers for zero-shot learning," in *CVPR*, 2016, pp. 5327–5336.

[87] S. Changpinyo, W.-L. Chao, and F. Sha, "Predicting visual exemplars of unseen classes for zero-shot learning," in *ICCV*, 2017, pp. 3496–3505.

[88] X. Wang, Y. Ye, and A. Gupta, "Zero-shot recognition via semantic embeddings and knowledge graphs," in *CVPR*, 2018, pp. 6857–6866.

[89] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *NeurIPS*, 2019.

[90] Y. Yu, Z. Ji, J. Han, and Z. Zhang, "Episode-based prototype generating network for zero-shot learning," in *CVPR*, 2020, pp. 14 032–14 041.

[91] S. Chen, Z. Hong, G.-S. Xie, W. Yang, Q. Peng, K. Wang, J. Zhao, and X. You, "Msdn: Mutually semantic distillation network for zero-shot learning," in *CVPR*, 2022.

[92] L. V. D. Maaten and G. E. Hinton, "Visualizing data using t-sne," *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008.

[93] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *NIPS*, 2013.

[94] W. Xu, Y. Xian, J. Wang, B. Schiele, and Z. Akata, "Vgse: Visually-grounded semantic embeddings for zero-shot learning," in *CVPR*, 2022.

[95] Z. Akata, S. E. Reed, D. Walter, H. Lee, and B. Schiele, "Evaluation of output embeddings for fine-grained image classification," in *CVPR*, 2015, pp. 2927–2936.

[96] M. Ye and Y. Guo, "Zero-shot classification with discriminative semantic representation learning," in *CVPR*, 2017, pp. 5103–5111.

[97] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *ICCV*, 2015.

[98] Y. Fu, T. Xiang, Y.-G. Jiang, X. Xue, L. Sigal, and S. Gong, "Recent advances in zero-shot recognition: Toward data-efficient understanding of visual content," *IEEE Signal Processing Magazine*, vol. 35, pp. 112–125, 2018.

**Shiming Chen** is currently a full-time Ph.D. student in the School of Electronic Information and Communications, Huazhong University of Sciences and Technology (HUST), China. His research results have expounded in prominent conferences and prestigious journals, such as NeurIPS, ICCV, CVPR, AAAI, IJCAI, IEEE TNNLS, IEEE TEVC, and etc. His current research interests span computer vision and machine learning with a series of topics, such as generative modeling and learning, zero-shot learning, and visual-and-language learning.

**Ziming Hong** is currently pursuing the M.Sc. degree in the School of Electronic Information and Communications(EIC), Huazhong University of Sciences and Technology(HUST), China. He received the B.E. degree in the School of Information Engineering, Wuhan University of Technology(WHUT), in 2019. His current research interests include graph learning and zero-shot learning.

**Wenjin Hou** is currently a full-time M.Sc. student at the School of Electronic Information and Communication, Huazhong University of Science and Technology (HUST), China. He received the B.E. degree in the School of Information Science and Engineering, Lanzhou University (LZU), in 2021. His research interests include zero-shot learning and generative modeling and learning in the field of computer vision.

**Guo-Sen Xie** is currently a Professor at the School of Computer Science and Engineering, Nanjing University of Science and Technology, China. He received his Ph.D. degree from the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China, in 2016. His research interests include computer vision and machine learning.

**Yibing Song** is with AI³ Institute, Fudan University. He was a senior researcher at Tencent AI Lab. He has obtained a Ph.D. degree from City University of Hong Kong, a MPhil degree from the same university, and a bachelor degree from University of Science and Technology of China. He has served as area chairs for CVPR, NeurIPS, and ICLR, served as reviewers for premier computer vision and machine learning conferences, and received multiple outstanding reviewer awards in CVPR 2018-2020, ECCV 2022, and NeurIPS 2019.

**Jian Zhao** is currently an Assistant Professor with the Institute of North Electronic Equipment, Beijing, China. He received his Ph.D. degree from the National University of Singapore (NUS) in 2019. He has served as the guest editor of PRL and Electronics, the presentation chair of the CICAI'21, the session chair of the ACM MM'21, and the invited reviewer of NSFC, T-PAMI, IJCV, NeurIPS, CVPR, etc. He has received the "2020-2022 Young Elite Scientist Sponsorship Program" from China Association for Science and Technology, and the "2021-2023 Beijing Young Elite Scientist Sponsorship Program" from Beijing Association for Science and Technology. His main research interests include deep learning, pattern recognition, computer vision and multimedia. He has published over 40 cutting-edge papers (e.g., T-PAMI, IJCV, NeurIPS, CVPR, etc.). He has received the nomination for the USERN Prize 2021, and won the Lee Hwee Kuan Award (Gold Award) on PREMIA'19 and the "Best Student Paper Award" on ACM MM'18 as the first author.
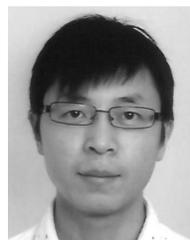
**Xinge You** (Senior Member, IEEE) is currently a Professor with the School of Electronic Information and Communications, Huazhong University of Science and Technology, Wuhan. He received the B.S. and M.S. degrees in mathematics from Hubei University, Wuhan, China, in 1990 and 2000, respectively, and the Ph.D. degree from the Department of Computer Science, Hong Kong Baptist University, Hong Kong, in 2004. His research results have expounded in 60+ publications at prestigious journals and prominent conferences, such as IEEE T-PAMI, T-IP, T-NNLS, NeurIPS, CVPR, ICCV, ECCV, and etc. He served/serves as an Associate Editor of the *IEEE Transactions on Cybernetics*, *IEEE Transactions on Systems, Man, Cybernetics:Systems*. His current research interests include image processing, wavelet analysis and its applications, pattern recognition, machine earning, and computer vision.

**Shuicheng Yan** (Fellow, IEEE) is currently the director of Sea AI Lab (SAIL) and group chief scientist of Sea. He is an Fellow of Academy of Engineering, Singapore, IEEE Fellow, ACM Fellow, IAPR Fellow. His research areas include computer vision, machine learning and multimedia analysis. Till now, he has published over 600 papers in top international journals and conferences, with Google Scholar Citation over 40,000 times and H-index 105. He had been among "Thomson Reuters Highly Cited Researchers" in 2014, 2015, 2016, 2018, 2019. Dr. Yan's team has received winner or honorable-mention prizes for 10 times of two core competitions, Pascal VOC and ImageNet (ILSVRC), which are deemed as "World Cup" in the computer vision community. Also his team won over 10 best paper or best student paper prizes and especially, a grand slam in ACM MM, the top conference in multimedia, including Best Paper Award, Best Student Paper Award and Best Demo Award.

**Ling Shao** (Fellow, IEEE) is the Chief Scientist of Terminus Group and the President of Terminus International. He was the founding CEO and Chief Scientist of the Inception Institute of Artificial Intelligence, Abu Dhabi, UAE. His research interests include computer vision, deep learning, medical imaging and vision and language. He is a fellow of the IEEE, the IAPR, the BCS and the IET.