

Hybrid Routing Transformer for Zero-Shot Learning

De Cheng[†], Gerong Wang[†], Bo Wang, Qiang Zhang[✉], Jungong Han, Dingwen Zhang.

Abstract—Zero-shot learning (ZSL) aims to learn models that can recognize unseen image semantics based on the training of data with seen semantics. Recent studies either leverage the global image features or mine discriminative local patch features to associate the extracted visual features to the semantic attributes. However, due to the lack of the necessary top-down guidance and semantic alignment for ensuring the model attending to the real attribute-correlation regions, these methods still encounter a significant semantic gap between the visual modality and the attribute modality, which makes their prediction on unseen semantics unreliable. To solve this problem, this paper establishes a novel transformer encoder-decoder model, called hybrid routing transformer (HRT). In HRT encoder, we embed an active attention, which is constructed by both the bottom-up and the top-down dynamic routing pathways to generate the attribute-aligned visual feature. While in HRT decoder, we use static routing to calculate the correlation among the attribute-aligned visual features, the corresponding attribute semantics, and the class attribute vectors to generate the final class label predictions. This design makes the presented transformer model a hybrid of 1) top-down and bottom-up attention pathways and 2) dynamic and static routing pathways. Comprehensive experiments on three widely-used benchmark datasets, namely CUB, SUN, and AWA2, are conducted. The obtained experimental results demonstrate the effectiveness of the proposed method.

Index Terms—Zero-Shot Learning, Hybrid Routing, Transformer, Attention.

I. INTRODUCTION

Deep learning has made great progress in a variety of vision tasks when the models are trained on large-scale labeled datasets. However, the real-world natural images follow a long-tailed distribution so that the data-hungry characteristic of CNN-based models limits their ability to recognize rare object classes, specially for the fine-grained animal species [1]. Meanwhile, an increasing number of newly defined visual concepts and products come to the fore so quickly, and the speed of data annotation for model training cannot keep up with the pace of new things emerging, thus the CNN-based models cannot be generalized to these new classes for testing. In contrast, a child can learn from only a few samples, summarize knowledge, and even draw inferences about other cases from one instance to recognize unseen objects. Therefore, building zero-shot learning (ZSL) models to transfer knowledge from seen classes to unseen classes is significant and indispensable.

Zero-shot learning (ZSL) mimics the human ability to recognize objects only from a description in terms of concepts in some semantic vocabulary [2], and aims to recognize the unseen classes, of which the labeled images are unavailable

during model training [3], [4]. Existing works on ZSL mainly leverage the global features [5] or patch features [6], [7] to construct visual-semantic alignment models. In these approaches, images and attributes are embedded with compatibility function. Despite good performances on coarse-grained datasets (e.g., Animal with Attribute dataset [3]), these approaches gradually degenerate when dealing with fine-grained datasets (e.g., Caltech-UCSD Birds-200-2011 dataset [8]), since much more local discriminative information is required to distinguish these fine-grained classes. Several recent works [9], [10], [11], [12], [13] try to focus on discriminative visual feature learning, by introducing attention mechanism into zero-shot classification problem, such as the spatial and channel attention [9], region attention [10], [14]. However, there still exists significant semantic gap between the visual modality and the attribute modality in the existing passive attention mechanisms as these methods generate attention weights purely in the bottom-up forward passing manner, which lacks the necessary top-down guidance and semantic alignment for attending to the real attribute-correlation regions. In this paper, We reveal a fundamental issue that establishing an active connection between the visual feature and attribute vector rather than a simple passive link is a key point to facilitate zero-shot learning.

Specifically, when trying to recognize an image from unseen classes, humans will involuntarily try to establish a connection between the attribute semantics with the corresponding local image regions. Besides, humans achieve semantic alignment by ruling out the irrelevant visual regions and locating the most relevant ones in a gradual way [11]. The above two phenomena motivate us to introduce the top-down guidance and dynamic routing connection into the attention mechanism. As shown in Figure 1 and bottom-left of Figure 2, the newly proposed active attention is constructed by both the bottom-up and the top-down connection pathways, and each pathway is formed by the dynamic routing rather than the conventionally used convolutional forward passing. Such an active attention works as the transformer encoder in our framework to learn the attribute-aligned visual features.

After obtaining the attribute-aligned visual features, we calculate the correlation between the attribute-aligned visual features and the corresponding attribute semantics, and generate the final class label predictions by multiplying the obtained correlation vector to the class attribute vectors (see bottom-right of Figure 2). When considering the attribute semantics as keys, the attribute-aligned visual features as queries, and the class attribute vectors as values, respectively, it is interesting to see that such a process can also be interpreted as a transformer decoding process that works in semantic space. As the input visual features of the transformer decoder have already been well aligned with the semantic space, the involved elements

De Cheng, Gerong Wang, Qiang Zhang, Xidian University, Xi'an, Shaanxi, P.R. China, Bo Wang is with the Tsinghua University, Beijing, P.R. China, Jungong Han is with the Aberystwyth University, Aberystwyth, UK, Dingwen Zhang is with the Northwestern Polytechnical University, Xi'an, Shaanxi, P.R. China.

[†]De Cheng and Gerong Wang are co-first authors.

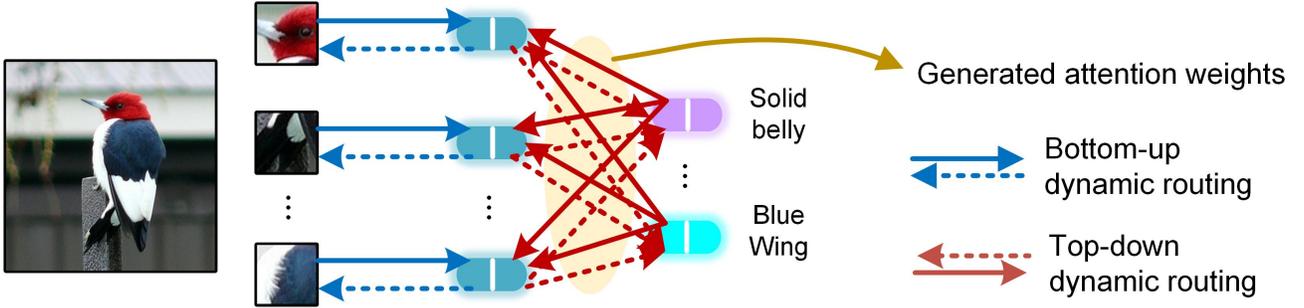


Fig. 1. Brief illustration of the propose active attention mechanism, which contains both the bottom-up dynamic routing process and the top-down dynamic routing process.

are connected with the static routing rather than the dynamic routing.

Based on the above-mentioned transformer encoder and decoder, we build a novel transformer model, called hybrid routing transformer (HRT). The overall learning model is shown in Figure 2. HRT mainly has three-fold novel properties: 1) From the perspective of the attention mechanism, it is not based on the commonly used passive self-attention. Instead, HRT is designed with the active semantic-guided attention. Particularly, the HRT encoder performs attention in both bottom-up and top-down manners. 2) From the perspective of the inner connection, HRT is formed by both the dynamic routing (in transformer encoder) and the static routing (in transformer decoder), which build different types of connection to the elements involved in the transformer. 3) From the perspective of zero-shot learning, we reveal an understudied yet important issue, i.e., the active connection between the visual feature and attribute vector, and build the first transformer-based zero-shot image recognition framework.

II. RELATED WORK

Zero-Shot Learning. Zero-shot learning (ZSL) aims to predict objects in unseen classes or both seen and unseen classes [15], the former is called traditional ZSL while the later is called generalized ZSL (GZSL). The core is to transfer knowledge learned from seen classes to unseen classes. The existing methods can be divided into three types: (1) Embedding methods [16], [5], [17], [18], [19], [6], [7], [20], [21], [22], [23], [24], which usually learn a compatibility function between image and class embedding spaces for similarity measurement. Specially, Akata *et al.* [25] propose a bilinear-style hinge loss to learn the compatibility function. Based on [25], Xian *et al.* [26] introduce non-linearity to ALE model. Following the structured SVM formulation, Akata *et al.* [5] designs a multiclass loss. Huynh *et al.* [27] leverage attribute semantic vectors to learn the association between images and attributes. (2) Generative methods [28], [29], [30], [31], [32], [33], [34], which aim to generate synthetic samples of unseen classes from semantic information and then set the ZSL problem as a supervised classification problem. Common generative methods use GAN [35], VAE [36] or flow-based generative models [37]. Xian *et al.* [38] directly generate image features conditioned on the class-level semantic descriptors. Felix *et*

al. [31] generate synthetic features by a multi-modal cycle-consistent GAN. (3) Gating Methods [39], [40], which use a gating based mechanism to separate the unseen samples from the seen samples for GZSL. Ideally, if the gate mechanism of binary classification is very effective, GZSL can be divided into a traditional ZSL problem to classify unseen samples and a supervised classification problem to classify seen samples. According to the experimental settings, the existing methods can be divided into two types: (1) Inductive ZSL, which only uses seen sample with labels in the training phase to classify unseen samples in the testing phase. (2) Transductive ZSL, which uses both seen samples with labels and unseen samples without labels in the training phase, which enables the model to use unseen visual features in the training phase to alleviate domain shift problem. Our experiments use the embedding method under the inductive ZSL setting.

Attention in ZSL. The aim of the attention mechanisms is to either highlight important local information or alleviate the influence of irrelevant and noisy information [11]. Ji *et al.* [11] weight different local features by a stacked attention mechanism, with access to the costly part annotations during training. Zhu *et al.* [9] weight different global features by learning multiple channel-wise attentions. Xie *et al.* [10] leverage attentive region embedding to learn the bilinear mapping to the semantic space. [27] is the closest competitor, which uses the passive attention to build similarity between features and class attribute vectors. However, these works either need part annotations or use passive semantic-unguided attention mechanism, thus enormous semantic gap still exists between the two unrelated modalities of image and attribute. In order to establish active correlation between image and attribute to capture discriminative features, we apply the capsule to transformer for ZSL, and build dynamic bottom-up and top-down attention mechanism by initializing high-level capsules with class-semantic vectors and performing low-level capsules with patch features. Our method proved to be very effective in subsequent experiments.

Capsule-Transformer. The transformer is proposed in [41] with the attention-based encoder-decoder architecture. It is successfully used in the natural language processing field [42] firstly and then extended to computer vision tasks [43], [44]. Capsule network was first introduced by [45], aiming to improve the ability of identifying spatial relationships and rotation of the CNN structure. Capsule is a group of neurons.

Using the dynamic routing method between two capsule layers, capsule network can match CNN in recognition results. Although, there appears few works [46], [47] to improve the transformer attention with capsule network for machine translation. The proposed HRT method has distinct properties with the existing models by using the active semantic-guided attention in both the bottom-up and top-down manner. It is also worth mention that this is the earliest work to establish a capsule-transformer-like framework for solving the ZSL problem.

III. HYBRID ROUTING TRANSFORMER

A. Problem Setting and Overall Framework

In zero-shot learning, we consider seen classes \mathcal{C}_s and unseen classes \mathcal{C}_u , where $\mathcal{C}_s \cap \mathcal{C}_u = \emptyset$. Specifically, we denote the training data as $\mathcal{D}^s = \{(x_i, y_i, \mathbf{z}_i)\}$, where x_i and y_i denote the training image and the corresponding label, and $\mathbf{z}_i = [z_1, \dots, z_A]$ represents the associated class attribute vectors. The target data contains images with both seen classes and unseen classes as well as the semantic vectors for each class.

In this work, given an input image, we first extract the basic feature representation for each image patch by ResNet-101. Then, we solve ZSL by a newly-designed transformer, where the transformer encoder extracts attribute-aligned visual features, while the transformer decoder projects the attribute-aligned visual features into classification scores. In particular, by designing bottom-up and top-down, dynamic and static routing layers in the transformer architecture, we build hybrid routing transformer to perform active interaction between images and attributes.

B. HRT Encoder

Given the patch features $\{\mathbf{f}^r\}_{r=1}^R$ of each input image, we equip HRT Encoder with attribute-guided dynamic routing layers to obtain the attribute-aligned visual features. Firstly, in order to obtain a compact description of each image patch, we cluster every patch feature into one capsule of a smaller dimension $d = 16$ by a bottom-up routing process, creating new patch capsule $\mathbf{g}^r \in \mathbb{R}^{1 \times d}$. In specific, given one patch feature vector $\mathbf{f}^r \in \mathbb{R}^{1 \times 2048}$, we convert it to primary capsules, each of which is represented by a 4×4 pose matrix, through a 1×1 learnable convolution on \mathbf{f}^r and obtain 128 primary capsules for each image region as illustrated in Figure 2. Then, we use a bottom-up routing process to facilitate the routing between the 128 primary (child) capsules and one parent capsule. The bottom-up routing is implemented based on the EM routing [48]. For child capsule i , firstly it is transformed by \mathbf{T}_{ij} to cast a vote $\mathbf{O}_{ij} = \mathbf{M}_i \mathbf{T}_{ij}$ for the parent capsule j . Then the non-linear routing process is performed by the EM algorithm, where the vectorized version of j^{th} pose matrix \mathbf{M}_j is the expectation of j^{th} Gaussian distribution. Denote $P_{i|j}$ as the probability density of the vectorized vote \mathbf{O}_{ij} under j^{th} Gaussian model. $P_{i|j}^h$ is h^{th} component of $P_{i|j}$, which

has variance $(\sigma_j^h)^2$ and mean μ_j^h . The h^{th} component of the probability density $P_{i|j}^h$ is computed as,

$$P_{i|j}^h = \frac{1}{\sqrt{2\pi} (\sigma_j^h)^2} \exp\left(-\frac{(\mathbf{O}_{ij}^h - \mu_j^h)^2}{2 (\sigma_j^h)^2}\right). \quad (1)$$

Then we calculate the activation of the parent capsule j based on the minimum description length principle :

$$\alpha_j = \text{logistic}\left(\lambda \left(\beta - \gamma \sum_i r_{ij} - \sum_h \text{cost}_j^h\right)\right), \quad (2)$$

$$\text{cost}_j^h = -\sum_i r_{ij} \ln P_{i|j}^h,$$

where β and γ are two learnable parameters, cost_j^h indicates the cost for activating the parent capsule j . $\sum_i r_{ij}$ calculates the amount of child capsules assigned to the parent capsule j . λ is a hyper-parameter. In EM routing, M-step and E-step run iteratively. The M-step computes the outputs of pose matrix and activation of the parent capsule j , and the E-step exports the possibility of child capsules assigned to the parent capsule j . We obtain the parent capsule \mathbf{g}^r as the compact representation for each image patch by iteratively calculating the pose matrix \mathbf{M} between the child and parent capsules in the EM routing process.

In order to attain active guidance of attribute semantics, we then use a top-down routing process to establish connections between the obtained image patch capsules and the global attribute semantic capsules transformed from the original attribute semantic vectors. Similar to [27], we apply GloVe [49] to extract the τ -dimensional attribute semantic vectors $\{\mathbf{v}_a\}_{a=1}^A$, which is followed by a demission reduction process based on the Factor Analysis [50]. Finally, we obtain the d -dimensional compact attribute semantic vectors $\{\tilde{\mathbf{v}}_a\}_{a=1}^A$, where $a \in [1, 2, \dots, A]$ is the attribute index. To build the top-down routing connection, we adopt the Inverted Dot-Product Attention routing process [51]. Such routing process first calculates the vote $\nu_{ij} = \mathbf{W}_{ij}^e \cdot \mathbf{p}_i$ for the child capsule \mathbf{p}_i . Then it computes the agreement as $\mathbf{o}_{ij} = \mathbf{p}_j^\top \cdot \nu_{ij}$ by the dot-product similarity between parent capsule \mathbf{p}_j and the vote. Then we update the parent capsules \mathbf{p}_j by:

$$\mathbf{p}_j = \text{LayerNorm}\left(\sum_i r_{ij} \nu_{ij}\right), r_{ij} = \frac{\exp(\mathbf{o}_{ij})}{\sum_{j'} \exp(\mathbf{o}_{ij'})}. \quad (3)$$

The above routing process will be performed several times to strength the agreement between the child and parent capsules.

It worth mentioning that, the parent capsules \mathbf{p}_j in our framework are initialized with $\{\tilde{\mathbf{v}}_a\}_{a=1}^A$, while other capsule networks perform zero or random initialization on them. The advantage of our approach is that it can realize an attribute-guided and active attention mechanism between the two-modality information. In fact, such a routing process simultaneously facilitates the semantic feature embedding and the semantic coherence computing. Under this circumstance, we treat the attribute semantic vector as attribute query \mathbf{Q} , the patch capsule as visual key \mathbf{K} . Then, the agreement between the child and parent capsules can properly reflect the

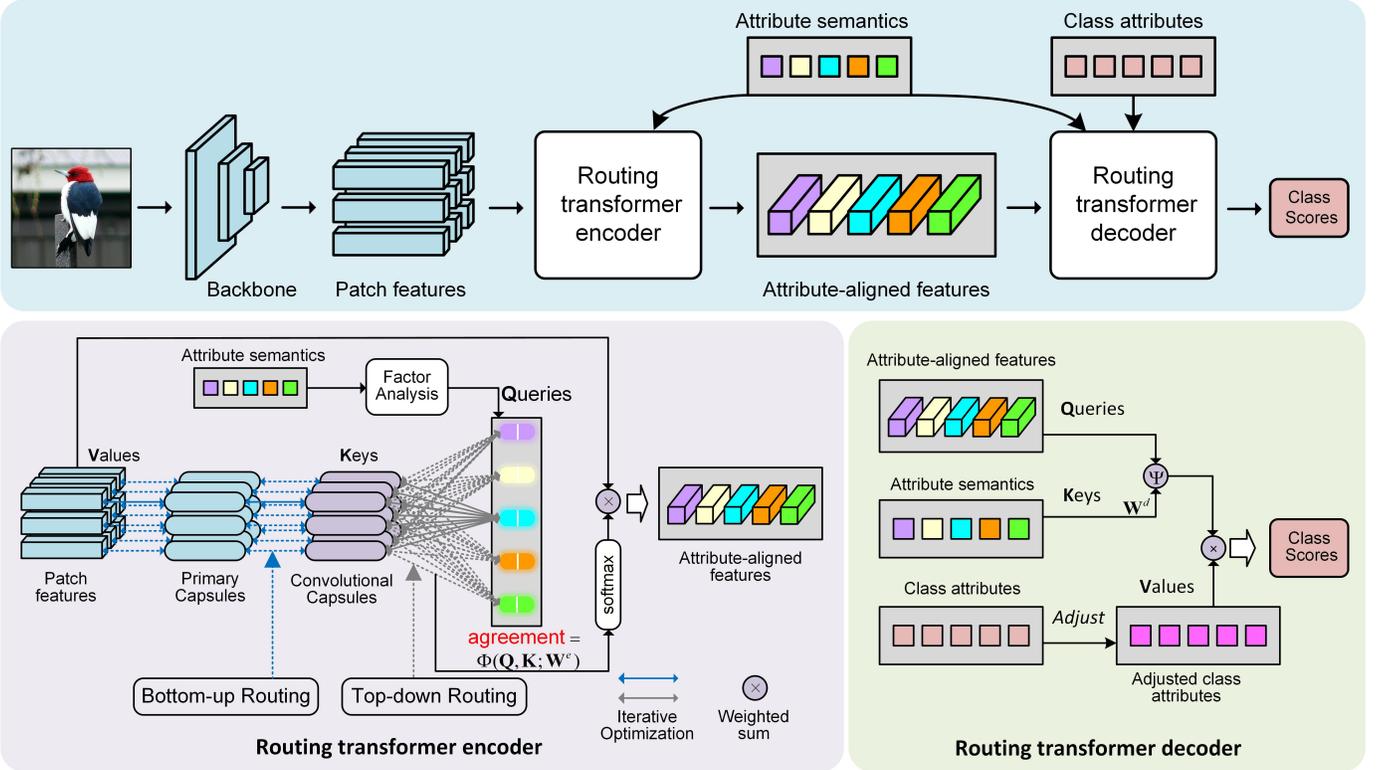


Fig. 2. Our Hybrid routing transformer consists of an *Semantic-guided dynamic transformer routing Part* as the encoder constructing the crucial semantic-aligned visual features of an image and an *Semantic-guided static transformer routing Part* as the decoder transferring the attribute-aligned features into classification scores under the guidance of class attribute vectors. Both of the two parts are trained jointly to establish the necessary top-down guidance and semantic alignment.

similarity or the relationship between \mathbf{Q} and \mathbf{K} . Since the agreement is implicitly parameterized by \mathbf{W}^e , we denote it as $\Phi(\mathbf{Q}, \mathbf{K}; \mathbf{W}^e) \in \mathbb{R}^{R \times A}$. Then, the final attribute-aligned visual features $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_1, \dots, \mathbf{h}_A] \in \mathbb{R}^{2048 \times A}$ can be obtained by:

$$\mathbf{H} = \text{HRT}^E(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\Phi(\mathbf{Q}, \mathbf{K}; \mathbf{W}^e))\mathbf{V}, \quad (4)$$

where $\mathbf{V} \in \mathbb{R}^{2048 \times R}$ indicates the value matrix, which is formed by the original R patch features $\{\mathbf{f}^r\}$.

C. HRT Decoder

Given the attribute-aligned visual features \mathbf{H} , we equip HRT Decoder with semantic-guided static routing layers to obtain the final class probability.

In the fine-grained recognition, there are so many attribute scores for classes while only a small portion of the attributes are crucial to distinguish different classes. In order to focus on the important attributes, we adjust the c -th class attribute vectors belonging by $\tilde{\mathbf{z}}^c = \text{sigmoid}(\Lambda^T \mathbf{W}_\beta \mathbf{H} \odot \mathbf{L}) \mathbf{z}^c$, where $\Lambda = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_A] \in \mathbb{R}^{\tau \times A}$, \mathbf{L} is the unitary matrix, \odot indicates the element-wise production. As the decoding process works in semantic space, we regard the attribute-aligned visual features in \mathbf{H} as the queries $\mathbf{Q} \in \mathbb{R}^{2048 \times A}$, the attribute semantic vectors $\{\mathbf{v}_a\}_{a=1}^A$ as the keys $\mathbf{K} \in \mathbb{R}^{\tau \times A}$, while the adjusted class attribute vectors $\{\tilde{\mathbf{z}}^c\}_{c=1}^C$ as the values $\mathbf{V} \in \mathbb{R}^{A \times C}$. Considering that the attribute-aligned visual features and the attribute semantic vectors have a

clear correspondence relationship, we define $\Psi(\mathbf{Q}, \mathbf{K}; \mathbf{W}^d) = \mathbf{I}_{1 \times A} \text{diag}(\mathbf{Q}^T \mathbf{W}^d \mathbf{K})$ to represent the content-aware attribute vectors, where $\mathbf{W}^d \in \mathbb{R}^{2048 \times \tau}$ is an embedding matrix between \mathbf{Q} and \mathbf{K} . Then, the final class scores of the input image can be obtained by measuring the coherence between the content-aware attribute vectors and the adjusted class attribute vectors:

$$\mathbf{s} = \text{HRT}^D(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \Psi(\mathbf{Q}, \mathbf{K}; \mathbf{W}^d)\mathbf{V}. \quad (5)$$

D. Training and Testing

For training the proposed deep model, we leverage three-fold loss function:

$$\mathcal{L}_{\text{HRT}} = \mathcal{L}_{ce} + \lambda_1 \mathcal{L}_{cal} + \lambda_2 \mathcal{L}_{reg}, \quad (6)$$

where λ_1 and λ_2 are hyper-parameters. The first loss \mathcal{L}_{ce} is the cross-entropy loss, which measures the consistency between the classification prediction and the ground-truth label. The second one is the calibration loss [27] $\mathcal{L}_{cal} = -\sum_c y^c \log p(s^c + \gamma_c)$, where $p(s^c) = \frac{\exp(s^c)}{\sum_{c' \in C} \exp(s^{c'})}$, s^c is the prediction of the c -th class and y^c is its corresponding label, γ_c is the hyper-parameter to balance the prediction score between the seen and unseen classes, which is different for seen and unseen classes. This loss aims to alleviate the biased prediction towards the seen categories [27]. The third one is the attribute regression loss $\mathcal{L}_{reg} = \|\varphi(\mathbf{Q}, \mathbf{K}; \mathbf{W}^d) - \mathbf{z}^{c*}\|_2^2$, where c^* indicates the ground-truth class index. This loss

constraints the predict attribute score vector to agree with the class attribute vector that is obtained by human statistic, which can help to improve the generalization capacity to the unseen classis.

In testing, we use $c^* = \operatorname{argmax}_{c \in \mathcal{C}}(s^c + \gamma_c)$ to obtain the predicted class belonging to either the seen classes or unseen ones.

IV. EXPERIMENTS

A. Setting

Datasets. We evaluate the performance of our model on three widely used ZSL benchmark datasets: *Caltech-UCSD Birds-200-2011* (CUB) [8], *SUN attributes* (SUN) [52] and *Animals with Attributes 2* (AWA2) [4]. The CUB [8] is a fine-grained dataset containing 11788 bird images from 150 seen classes and 50 unseen classes with 312 expert-annotated attributes. Although it contains discriminative attribute location annotation to distinguish fine-grained classes, our model works under the weakly supervised setting without the part annotations unlike [11]. SUN [52] is another fine-grained dataset containing 14204 scene images, from 645 seen classes and 72 unseen classes with 102 annotated attributes. Different from the above two datasets, AWA2 [4] dataset is in the coarse-grained level consisting of 37322 animal images, from 40 seen classes and 10 unseen classes with 85 attributes.

Evaluation Metrics. In this work, we measure the average per-class top-1 accuracy of our method on both traditional ZSL and GZSL setting. For traditional ZSL task, we evaluate test images only from unseen classes with T1 (top-1 accuracy). For GZSL task, we evaluate test images from both seen classes and unseen classes. Following the protocol proposed in [4], we report tr and ts as the average per-class top-1 accuracy of test images from seen classes and unseen classes, respectively. H is computed as the hamonic mean of tr and ts to measure the comprehensive performance, which can be calculated by $H \triangleq 2 \times \frac{tr \times ts}{tr + ts}$.

Visual features. We obtain the patch features at the last convolution layer of ResNet-101 [53] model pre-trained on the ImageNet-1K [54] dataset. For CUB, it contains so many similar sub-categories under a big category of birds which needs much abundant visual information to classify them. Therefore, in order to obtain richer visual features, we modify the stride of conv 5_x layer in ResNet-101 from 2 to 1, and get the patch features of size $14 \times 14 \times 2048$. The differences among the categories in SUN and AWA2 are relatively large, finer features may confuse the correspondence between visual features and attributes, and eventually deteriorate classification accuracy. Thus, we only adjusted the stride of feature extraction process for CUB dataset.

Implementation Details. Following [4], we adopt ResNet-101 [53] pre-trained on ImageNet-1K [54] as the backbone for feature extraction without fine-tuning. Given the input image of size 224×224 , we will obtain patch features of size $7 \times 7 \times 2048$ or $14 \times 14 \times 2048$ for different datasets. We use RMSprop [55] optimization method by setting momentum as 0.9, weight decay of 10^{-4} and the initial learning rate of 10^{-3} . The coefficients λ_1 and λ_2 in Eq.(6) are set as 0.1 and 0.033,

respectively. The factor γ_c in calibration loss for seen classes is different, which is -0.5 for CUB and SUN dataset, -0.8 for AWA2 dataset, while $\gamma_c = 1$ for unseen classes on all the three datasets. The model is implemented based on PyTorch platform [56], training on a single 2080 Ti GPU card.

B. Comparison with the state-of-the-art

In the experiment, we compare our proposed HRT method with several state-of-the-art embedding methods on both ZSL and GZSL settings. We report the top-1 accuracy and harmonic mean of each method in Table I, where “-” indicates that the results are not reported. These methods, CONSE [57], DEVISE [58], ALE [25], SJE [5], ESZSL [59], SSE [24], SYNC [60], LATEM [26], SAE [61], learn a compatibility function between visual image and attributes for similarity measurement. Meanwhile, these methods, SGMA [9], LFGAA [62], AREN [10] and DAZLE [27], introduce the attention mechanism to embedding methods. D-VAE [63] and GCM-CF [64] are generative models. Extensive experiments demonstrate the effectiveness of our proposed method, and the experiment results on all the three datasets show that our proposed method, achieve superior performances to the existing state-of-the-art methods in most cases.

Generalized Zero-shot Learning Results. As shown in Table I, the proposed HRT method acheives superior results on AWA2, CUB and SUN datasets for generalized ZSL task in most cases. We gain impressive results for the harmonic mean (H), where 62.8% on CUB, and 67.4% on AWA2. On SUN dataset, we obtains 53.2% for top-1 accuracy on unseen classes, which is also much better than other methods. In particular, our method significantly outperforms other algorithms on unseen accuracy, surpassing the state-of-the-art model on unseen accuracy by 4.4% on CUB dataset. However, since our capsule attention model is a kind of the dense attention model, lacking the training samples makes the proposed method cannot achieve the best harmonic mean on SUN dataset, where this dataset contains only 16 training samples for each seen classes. Furthermore, DAZLE [27] is the closest competitor, which uses the passive attention to build similarity. It causes confused relationship and cannot bridge the semantic gap between the visual image and the attribute descriptions. Different from it, we construct an active semantic-guided attention mechanism which helps high-level semantic actively focus on relevant visual features. As a result, we achieve superior results, improving harmonic mean (H) of 4.7% on CUB and 2.5% on SUN dataset absolutely, compared to DAZLE [27]. Notice that LFGAA* [62] achieves the best seen accuracy, while our HRT model obtains the best unseen accuracy with relatively high seen accuracy. Moreover, the experiment results in Table I show that our proposed HRT model can balance the performance of both seen and unseen classes greatly, compared with these methods, AREN* [10], LFGAA* [62], SGMA* [9], LATEM [26]. This maybe due to that our proposed HRT method constructs an active attention

¹SEM of HRT means the Standard Error of mean accuracy of our model under 5 random seed experiments.

TABLE I

COMPARISON OF OUR METHOD WITH THE STATE-OF-THE-ART EMBEDDING METHODS ON CUB, SUN AND AWA2. THE METHODS WITH * DENOTE FINE-TUNING THE BACKBONE WEIGHTS, OTHERWISE FIXING THEM. WE MEASURE TOP-1 ACCURACY(T1) IN ZSL SETTING, AND TOP-1 ACCURACY ON SEEN/UNSEEN (TR/TS) CLASSES AND THEIR HARMONIC MEAN (H) IN GZSL SETTING. BOLD FONT DENOTE THE BEST RESULTS.

Method	Zero-shot Learning			Generalized Zero-shot Learning								
	AWA2	CUB	SUN	AWA2			CUB			SUN		
	T1	T1	T1	tr	ts	H	tr	ts	H	tr	ts	H
CONSE [57]	44.5	34.3	38.8	90.6	0.5	1.0	72.2	1.6	3.1	39.9	6.8	11.6
DEVISE [58]	59.7	52	56.5	74.7	17.1	27.8	53.0	23.8	32.8	27.4	16.9	20.9
ALE [25]	62.5	54.9	58.1	81.8	14	23.9	62.8	23.7	34.4	33.1	21.8	26.3
SJE [5]	61.9	53.9	53.7	73.9	8.0	14.4	59.2	23.5	33.6	30.5	14.7	19.8
ESZSL [59]	58.6	53.9	54.5	77.8	5.9	11.0	63.8	12.6	21	27.9	11.0	15.8
SSE [24]	61.0	43.9	51.5	82.5	8.1	14.8	46.9	8.5	14.4	36.4	2.1	4.0
SYNC [60]	46.6	55.6	56.3	90.5	10.0	18.0	70.9	11.5	19.8	43.3	7.9	13.4
LATEM [26]	55.8	49.3	55.3	77.3	11.5	20.0	57.3	15.2	24	28.8	14.7	19.5
SAE [61]	54.1	33.3	40.3	82.2	1.1	2.2	54.0	7.8	13.6	18.0	8.8	11.8
D-VAE [63]	-	-	-	80.2	56.9	66.6	58.2	51.1	54.4	47.6	36.6	41.4
GCM-CF [64]	-	-	-	75.1	60.4	67.0	59.7	61.0	60.3	37.8	47.9	42.2
SGMA* [9]	68.8	71.0	-	87.1	37.6	52.5	71.3	36.7	48.5	-	-	-
LFGAA* [62]	68.1	67.6	61.5	93.4	27.0	41.9	80.9	36.2	50.0	40.4	18.5	25.3
AREN* [10]	67.9	71.8	60.6	92.9	15.6	26.7	78.7	38.9	52.1	38.8	19.0	25.5
DAZLE [27]	67.9	65.9	59.3	75.7	60.3	67.1	59.6	56.7	58.1	24.3	52.3	33.2
HRT(ours)	67.3	71.7	63.9	78.7	58.9	67.4	63.5	62.1	62.8	26.9	53.2	35.7
SEM of HRT ¹	±0	±0.01	±0.04	±0.01	±0.02	±0.02	±0.01	±0.08	±0.03	±0.05	±0.01	±0.04

mechanism and facilitates a more effective knowledge transfer from seen classes to unseen classes.

Zero-shot Learning Results. In ZSL, training classes are disjoint with testing classes. As shown in Table I, our proposed method presents a significant improvement compared with the state-of-the-art methods in most cases, i.e., 2.4% on SUN. The SUN dataset is a scene dataset, our method is able to understand the abstract content of scenes more deeply, through constructing the hybrid dynamic top-down and bottom-up attention pathways between visual image and the attribute. On CUB and AWA2 datasets, we gain comparable performance with the best method AREN* [10] and SGMA* [9].

C. Ablations

To illustrate the effectiveness of the proposed capsule network based HRT encoder framework, we make detailed performance comparison between the traditional CNN based HRT encoder and the capsule network based encoder framework on CUB, SUN and AWA2 datasets as shown in Table II. We denote them as HRT-CNNNet and HRT-CapsuleNet, respectively. The methods HRT-CNNNet and HRT-CapsuleNet are trained only with cross-entropy and calibration loss, while HRT-Ours is trained by the whole loss in Eq 6. We can clearly see that the capsule network based encoder framework is better than that of the traditional CNN based HRT encoder framework, HRT-CapsuleNet outperforms the baseline HRT-CNNNet by 2.1% under the H(harmonic mean) measurement in the GZSL setting, while 2.4% under the T1 measurement in the ZSL setting.

To further measure the influences of different components in HRT, we perform ablation study on CUB dataset under both ZSL and GZSL settings. As shown in Table III, the first line means that we train HRT only with cross-entropy loss and the patch features in hybrid routing transformer encoder

and decoder are performed without any attention mechanism. By adding hybrid routing transformer encoder, hybrid routing transformer decoder, \mathcal{L}_{cal} loss and \mathcal{L}_{reg} loss gradually, the GZSL results demonstrate that these four components improve accuracy by a huge margin over baseline, i.e., 3.8% (tr), 60.2% (ts), 59.1% (H) and 20.3% (T1), absolutely. The second line means that we perform the semantic-guided hybrid routing transformer in encoder, and the harmonic mean is significantly improved by 13.1%, which is achieved by constructing both the bottom-up and the top-down dynamic routing pathways to generate the aligned features. The third line represents that by adding the semantic-guided static transformer routing in decoder, and the model further improves the performances. In addition, the calibration loss is an essential part to boost accuracy under GZSL setting. Finally, through attribute regression loss, we further strengthen supervision and correct the classification results, boosting the harmonic mean by 3.4%. The above four components constitute the model together, improving the performance under both ZSL and GZSL settings.

D. Hyper-Parameter Selection

For the purpose of establishing the necessary top-down guidance and semantic alignment for attending to the real attribute-correlation regions, we build the HRT model via bottom-up and top-down routings. To observe the influence of the iteration times of these two routings, we conduct experiments with various iteration values, i.e. {1, 2, 3, 4, 5}. General ZSL results on CUB dataset are illustrated in Figure 3. In Figure 3 (a), we set the EM routing iteration time to 1, and the top-down dynamic attention routing iteration time changes from 1 to 5. We get the best results with the top-down dynamic attention routing iteration of 2. In Figure 3 (b), we fix the top-down dynamic attention routing iteration time to 2 and study the performances under different EM routing

TABLE II

PERFORMANCE ANALYSIS OF OUR HRT FRAMEWORK WITH CNN AND THE CAPSULE NETWORK IN THE TRANSFORMER ENCODER ON CUB, SUN AND AWA2, WE DENOTE THEM AS HRT-CNNNET AND HRT-CAPSULENET RESPECTIVELY. THE METHODS HRT-CNNNET AND HRT-CAPSULENET ARE TRAINED ONLY WITH CROSS-ENTROPY AND CALIBRATION LOSS, WHILE HRT-OURS IS TRAINED BY THE WHOLE LOSS IN EQ 6. WE MEASURE TOP-1 ACCURACY(T1) IN ZSL SETTING, AND TOP-1 ACCURACY ON SEEN/UNSEEN (TR/TS) CLASSES AND THEIR HARMONIC MEAN (H) IN GZSL SETTING.

Method	Zero-shot Learning			Generalized Zero-shot Learning								
	AWA2	CUB	SUN	AWA2			CUB			SUN		
	T1	T1	T1	tr	ts	H	tr	ts	H	tr	ts	H
HRT-CNNNet	66.9	64.3	59.0	75.9	59.3	66.6	59.4	54.8	57.0	23.6	51.8	32.4
HRT-CapsuleNet	67.7	69.3	60.4	75.3	60.2	66.9	62.4	60.0	61.2	25.8	50.8	34.2
HRT-Ours	67.7	71.7	63.9	78.7	58.9	67.4	63.5	62.1	62.8	26.9	53.2	35.7

TABLE III

ABLATION RESULTS FOR (GENERALIZED) ZERO-SHOT LEARNING ON CUB DATASET UNDER THE SAME BASELINES.

Hybrid Routing Transformer Encoder	Hybrid Routing Transformer Decoder	\mathcal{L}_{cal}	\mathcal{L}_{reg}	tr	ts	H	T1
✗	✗	✗	✗	59.7	1.9	3.7	51.4
✓	✗	✗	✗	66.4	9.6	16.8	63.1
✓	✓	✗	✗	69.6	10.0	17.5	67.7
✓	✓	✓	✗	62.4	60.0	61.2	69.8
✓	✓	✓	✓	63.5	62.1	62.8	71.7

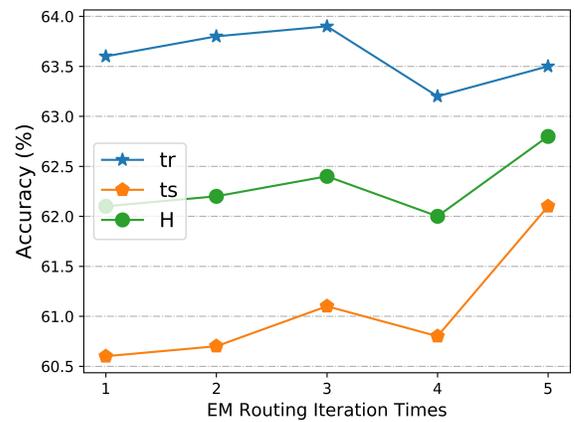
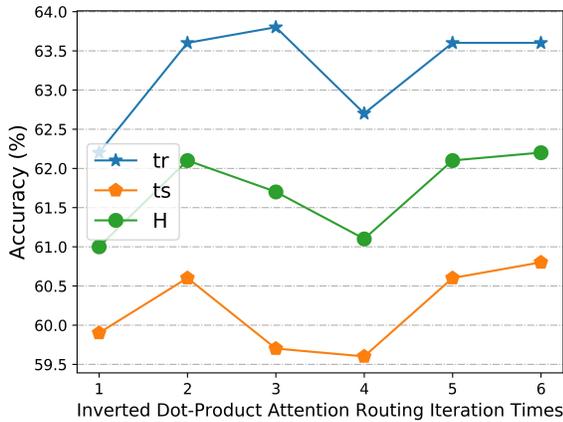


Fig. 3. (a) shows the changes of tr, ts and H with different Inverted Dot-Product Attention routing iteration times under CUB dataset while EM routing iteration time is equal to 1. (b) displays the result with variational EM routing iteration times while Inverted Dot-Product Attention routing iteration time is set to 2.

iteration times. We obtain the best results with EM routing iteration of 5. The above experiment results demonstrate the superiority of the proposed dynamic bottom-up and top-down routing mechanism, which helps HRT model establish an active connection between the visual feature and attribute vector rather than a simple passive mapping.

E. Qualitative results

We compared the feature visualization results of the proposed HRT model with the baseline passive attention model DAZLE[27], by contacting image feature and attribute on the CUB dataset. Figure 4 visualizes the agreement maps for examples from unseen and seen classes. On the whole, the attention maps generated by the HRT method are more concentrated and contacting image feature and attribute more accurately. For the bird “Evening Grosbeak” of unseen classes,

DAZLE tends to generate inaccurate or deflected attention maps covering the whole bird. For example, with the DAZLE model, the attributes “Bill Color Yellow”, “Throat Color Yellow”, “Belly Color Buff” and “Leg Color Pink” concentrate on the false visual information, and the attribute “Wing Pattern Striped” doesn’t capture relevant visual information. For the bird “Black Throated Blue Warbler” of the seen classes, though DAZLE can generate attention maps utilizing the learned knowledge, it can’t capture the true attribute-aligned features and its attention maps are more dispersive. On the other hand, the attention maps of our HRT model focuses on the attribute-aligned visual features more accurately. These examples demonstrate that our hybrid routing transformer model constructs an active and effective connection between image feature and attribute, and thus alleviate the semantic gap effectively between two different modality information for zero-shot learning.

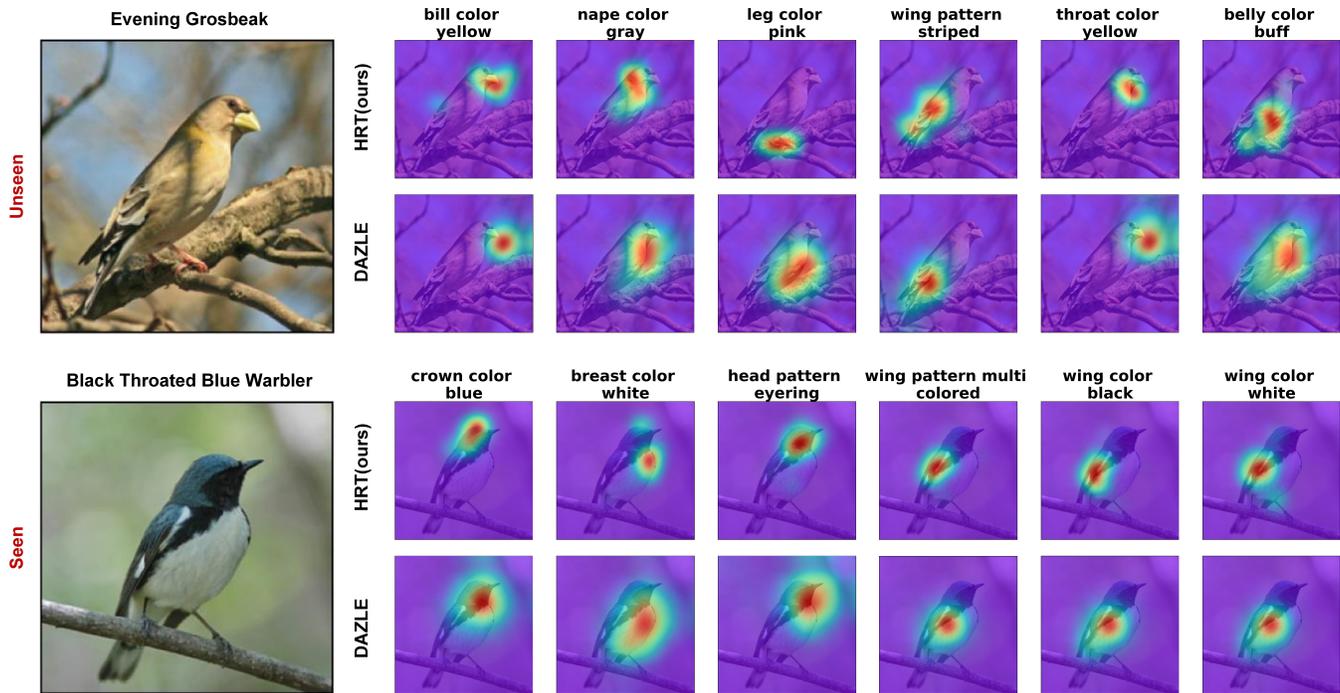


Fig. 4. Visualization of attribute-aligned feature locating for examples from unseen and seen classes on CUB dataset.

V. CONCLUSIONS

In this work, we propose a hybrid routing transformer (HRT) framework for ZSL and GZSL tasks. We are the first to apply the routing-based transformer for ZSL, by constructing semantic-guided active mechanism to alleviate the semantic gap between the visual modality and attribute modality effectively, and facilitate knowledge transfer. The proposed HRT model is a novel encoder-decoder framework. In the HRT encoder part, we utilize both the bottom-up and top-down dynamic routing pathways to generate attribute-aligned visual features. While in the HRT decoder part, we take the semantic-guided static routing to transfer attribute-aligned features into classification scores under the guidance of class attribute vectors. Extensive experiments suggest that the proposed active bottom-up and top-down dynamic routing pathway can help improve the transformer, thus we can extend our method to other research fields in the future.

REFERENCES

- [1] Z. Han, Z. Fu, and J. Yang, "Learning the redundancy-free features for generalized zero-shot object recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 865–12 874.
- [2] P. Morgado and N. Vasconcelos, "Semantically consistent regularization for zero-shot recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6060–6069.
- [3] C. H. Lampert, H. Nickisch, and S. Harmeling, "Attribute-based classification for zero-shot visual object categorization," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 3, pp. 453–465, 2013.
- [4] Y. Xian, B. Schiele, and Z. Akata, "Zero-shot learning-the good, the bad and the ugly," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4582–4591.
- [5] Z. Akata, S. Reed, D. Walter, H. Lee, and B. Schiele, "Evaluation of output embeddings for fine-grained image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2927–2936.
- [6] G.-S. Xie, L. Liu, F. Zhu, F. Zhao, Z. Zhang, Y. Yao, J. Qin, and L. Shao, "Region graph embedding network for zero-shot learning," in *European Conference on Computer Vision*. Springer, 2020, pp. 562–580.
- [7] W. Xu, Y. Xian, J. Wang, B. Schiele, and Z. Akata, "Attribute prototype network for zero-shot learning," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [8] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona, "Caltech-ucsd birds 200. california institute of technology," CNS-TR-2010-001, Tech. Rep., 2010.
- [9] Y. Zhu, J. Xie, Z. Tang, X. Peng, and A. Elgammal, "Semantic-guided multi-attention localization for zero-shot learning," in *Advances in Neural Information Processing Systems*, 2019, pp. 14 943–14 953.
- [10] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9384–9393.
- [11] Z. Ji, Y. Fu, J. Guo, Y. Pang, Z. M. Zhang *et al.*, "Stacked semantics-guided attention model for fine-grained zero-shot learning," in *Advances in Neural Information Processing Systems*, 2018, pp. 5995–6004.
- [12] N. L. J. R. D. Z. Peiliang Huang, Junwei Han, "Scribble-supervised video object segmentation," *IEEE/CAA Journal of Automatica Sinica*, p. 10.1109/JAS.2021.1004210, 2021.
- [13] G. C. M.-H. Y. Dingwen Zhang, Junwei Han, "Weakly supervised object localization and detection: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, p. 10.1109/TPAMI.2021.3074313, 2021.
- [14] J. Y. J. H. Dingwen Zhang, Wenyuan Zeng, "Weakly supervised object detection using proposal- and semantic-level relationships," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, p. 10.1109/TPAMI.2021.3074313, 2020.
- [15] Y. Xian, S. Sharma, B. Schiele, and Z. Akata, "f-vaegan-d2: A feature generating framework for any-shot learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 275–10 284.
- [16] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid, "Label-embedding for image classification," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 7, pp. 1425–1438, 2015.
- [17] Y. Annadani and S. Biswas, "Preserving semantic relations for zero-shot

- learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7603–7612.
- [18] S. Liu, J. Chen, L. Pan, C.-W. Ngo, T.-S. Chua, and Y.-G. Jiang, “Hyperbolic visual embedding learning for zero-shot recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9273–9281.
- [19] Y. Liu, Q. Gao, J. Li, J. Han, and L. Shao, “Zero shot learning via low-rank embedded semantic autoencoder,” in *IJCAI*, 2018, pp. 2490–2496.
- [20] X. Xu, F. Shen, Y. Yang, D. Zhang, H. Tao Shen, and J. Song, “Matrix tri-factorization with manifold regularizations for zero-shot learning,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3798–3807.
- [21] S. M. W. L. G. K. J. Y. Chen Gong, Dacheng Tao, “Multi-modal curriculum learning for semi-supervised image classification,” *IEEE Transactions on Image Processing (TIP)*, pp. 3249–3260, 2016.
- [22] F. Zhang and G. Shi, “Co-representation network for generalized zero-shot learning,” in *International Conference on Machine Learning*, 2019, pp. 7434–7443.
- [23] H. Zhang and P. Koniusz, “Zero-shot kernel learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7670–7679.
- [24] Z. Zhang and V. Saligrama, “Zero-shot learning via semantic similarity embedding,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4166–4174.
- [25] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid, “Label-embedding for attribute-based classification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 819–826.
- [26] Y. Xian, Z. Akata, G. Sharma, Q. Nguyen, M. Hein, and B. Schiele, “Latent embeddings for zero-shot classification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 69–77.
- [27] D. Huynh and E. Elhamifar, “Fine-grained generalized zero-shot learning via dense attribute-based attention,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4483–4493.
- [28] Y. Yu, Z. Ji, J. Han, and Z. Zhang, “Episode-based prototype generating network for zero-shot learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 035–14 044.
- [29] J. Wu, T. Zhang, Z.-J. Zha, J. Luo, Y. Zhang, and F. Wu, “Self-supervised domain-aware generative network for generalized zero-shot learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 767–12 776.
- [30] L. Chen, H. Zhang, J. Xiao, W. Liu, and S.-F. Chang, “Zero-shot visual recognition using semantics-preserving adversarial embedding networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1043–1052.
- [31] R. Felix, V. B. Kumar, I. Reid, and G. Carneiro, “Multi-modal cycle-consistent generalized zero-shot learning,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 21–37.
- [32] J. Y. M. S. Chen Gong, Jian Yang, “Centroid estimation with guaranteed efficiency: A general framework for weakly supervised learning,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.
- [33] V. Kumar Verma, G. Arora, A. Mishra, and P. Rai, “Generalized zero-shot learning via synthesized examples,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4281–4289.
- [34] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” *Advances in neural information processing systems*, vol. 26, pp. 3111–3119, 2013.
- [35] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” *arXiv preprint arXiv:1406.2661*, 2014.
- [36] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [37] L. Dinh, D. Krueger, and Y. Bengio, “Nice: Non-linear independent components estimation,” *arXiv preprint arXiv:1410.8516*, 2014.
- [38] Y. Xian, T. Lorenz, B. Schiele, and Z. Akata, “Feature generating networks for zero-shot learning,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5542–5551.
- [39] Y. Atzmon and G. Chechik, “Adaptive confidence smoothing for generalized zero-shot learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 671–11 680.
- [40] X. Chen, X. Lan, F. Sun, and N. Zheng, “A boundary based out-of-distribution classifier for generalized zero-shot learning,” in *European Conference on Computer Vision*. Springer, 2020, pp. 572–588.
- [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *Proceedings of the 31th Conference on Neural Information Processing Systems*, pp. 5998–6008, 2017.
- [42] K. L. K. T. Jacob Devlin, Ming-Wei Chang, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2019.
- [43] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [44] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *European Conference on Computer Vision*. Springer, 2020, pp. 213–229.
- [45] G. S. N. U. A. K. S. Z. Nicolas Carion, Francisco Massa, “End-to-end object detection with transformers,” *arXiv:2005.12872*, 2020.
- [46] M. U. E. S. Kehai Chen, Rui Wang and T. Zhao, “Syntax directed attention for neural machine translation,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [47] H. Z. Sufeng Duan1, Juncheng Cao, “Capsule-transformer for neural machine translation,” *arXiv:2004.14649v1*, 2020.
- [48] G. E. Hinton, S. Sabour, and N. Frosst, “Matrix capsules with em routing,” in *International conference on learning representations*, 2018.
- [49] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [50] J. C. Loehlin, *Latent variable models*. hillsdale, nj: erlbaum, 1987.
- [51] Y.-H. H. Tsai, N. Srivastava, H. Goh, and R. Salakhutdinov, “Capsules with inverted dot-product attention routing,” *arXiv preprint arXiv:2002.04764*, 2020.
- [52] G. Patterson, C. Xu, H. Su, and J. Hays, “The sun attribute database: Beyond categories for deeper scene understanding,” *International Journal of Computer Vision*, vol. 108, no. 1-2, pp. 59–81, 2014.
- [53] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [54] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [55] T. Tieleman and G. Hinton, “Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude,” *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.
- [56] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *arXiv preprint arXiv:1912.01703*, 2019.
- [57] 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,” *arXiv preprint arXiv:1312.5650*, 2013.
- [58] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, “Devise: A deep visual-semantic embedding model,” in *Advances in neural information processing systems*, 2013, pp. 2121–2129.
- [59] B. Romera-Paredes and P. Torr, “An embarrassingly simple approach to zero-shot learning,” in *International Conference on Machine Learning*, 2015, pp. 2152–2161.
- [60] S. Changpinyo, W.-L. Chao, B. Gong, and F. Sha, “Synthesized classifiers for zero-shot learning,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5327–5336.
- [61] E. Kodirov, T. Xiang, and S. Gong, “Semantic autoencoder for zero-shot learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3174–3183.
- [62] Y. Liu, J. Guo, D. Cai, and X. He, “Attribute attention for semantic disambiguation in zero-shot learning,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6698–6707.
- [63] X. Li, Z. Xu, K. Wei, and C. Deng, “Generalized zero-shot learning via disentangled representation,” 2021.
- [64] Z. Yue, T. Wang, H. Zhang, Q. Sun, and X.-S. Hua, “Counterfactual zero-shot and open-set visual recognition,” *arXiv preprint arXiv:2103.00887*, 2021.