

Fine-Grained Representation Learning and Recognition by Exploiting Hierarchical Semantic Embedding

Tianshui Chen
Sun Yat-sen University
tianshuichen@gmail.com

Wenxi Wu
Sun Yat-sen University
ngmanhei@foxmail.com

Yuefang Gao
South China Agricultural University
gaoyuefang@scau.edu.cn

Le Dong
University of Electronic Science and
Technology of China
ledong@uestc.edu.cn

Xiaonan Luo
Guilin University of Electronic
Technology
luoxn@guet.edu.cn

Liang Lin*
Sun Yat-sen University
linliang@ieee.org

ABSTRACT

Object categories inherently form a hierarchy with different levels of concept abstraction, especially for fine-grained categories. For example, birds (*Aves*) can be categorized according to a four-level hierarchy of order, family, genus, and species. This hierarchy encodes rich correlations among various categories across different levels, which can effectively regularize the semantic space and thus make prediction less ambiguous. However, previous studies of fine-grained image recognition primarily focus on categories of one certain level and usually overlook this correlation information. In this work, we investigate simultaneously predicting categories of different levels in the hierarchy and integrating this structured correlation information into the deep neural network by developing a novel Hierarchical Semantic Embedding (HSE) framework. Specifically, the HSE framework sequentially predicts the category score vector of each level in the hierarchy, from highest to lowest. At each level, it incorporates the predicted score vector of the higher level as prior knowledge to learn finer-grained feature representation. During training, the predicted score vector of the higher level is also employed to regularize label prediction by using it as soft targets of corresponding sub-categories. To evaluate the proposed framework, we organize the 200 bird species of the Caltech-UCSD birds dataset with the four-level category hierarchy and construct a large-scale butterfly dataset that also covers four level categories. Extensive experiments on these two and the newly-released VegFru datasets demonstrate the superiority of our HSE framework over the baseline methods and existing competitors.

KEYWORDS

Semantic Embedding, Fine-Grained Image Recognition, Category Hierarchy

*Tianshui Chen and Wenxi Wu contribute equally to this work and share first-authorship. Corresponding author is Liang Lin.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '18, Seoul, Republic of Korea

© 2018 ACM. 978-1-4503-5665-7/18/10...\$15.00

DOI: 10.1145/3240508.3240523

1 INTRODUCTION

Object categories inherently form a hierarchy with different levels of concept abstraction, in which nodes closer to the root of the hierarchy refer to more abstract concepts while nodes closer to the leaves refer to finer-grained concepts. This hierarchy organization is especially important and obvious for fine-grained categories. For example, the fine-grained categories of birds (*Aves*) can be organized with a four-level hierarchy of order, family, genus and species, where an order consists of several families while a family consists of several genera, and so on. This category hierarchy provides very rich semantic correlations among categories across different levels, which can effectively regularize semantic space and provide extra guidance to attend more subtle regions for better recognition. For example, to recognize the fine-grained category of a given object (e.g., the species of a bird), we might first recognize its superclass (e.g., genus). Then, we prefer to concentrate on the fine-grained categories that are subject to this superclass and fixate on object parts that are more distinguishable among these fine-grained categories.

Existing methods on fine-grained image recognition (FGIR) primarily focus on classifying categories of one particular level, e.g., categorizing 200 species of birds [25, 50] or 431 models of cars [15], and usually overlook this correlation information. In this work, we simultaneously predict categories of all levels in the hierarchy, and integrate this structured correlation information into the deep neural network to progressively regularize label prediction and guide representation learning. To this, we formulate a novel Hierarchical Semantic Embedding (HSE) framework that orderly predicts the score vector of each level, from highest to lowest. At each level, it incorporates the predicted score vector of the higher level as prior knowledge to learn finer-grained feature representation. This is implemented by a semantic guided attentional mechanism that learns to fixate on more discriminative regions for better distinguishing. During training, we also utilize the predicted score vector of the higher level as soft targets to regularize the label prediction, thus that the predicted result at this level finely accords with that predicted at the higher level.

Caltech-UCSD birds dataset [39] is the most widely used benchmark for evaluating the FGIR task. To evaluate our proposed HSE framework on this benchmark, we organize the 200 bird categories with a four-level hierarchy of 13 orders, 37 families, 122 genera, and 200 species according to the ornithological systematics [32, 33]. In addition, we also create a new large-scale butterfly (namely Butterfly-200) dataset that also covers four-level categories for

multi-granularity image recognition. Currently, this dataset consists of 200 prevalent species of butterflies, which are grouped into 116 genera, 23 sub-families, and 5 families according to the insect taxonomy [34, 38]. It contains 25,279 images in total and at least 30 images per species. It’s worth noting that these category hierarchies can be obtained from the literature of taxonomy [33, 38] or directly retrieved from Wikipedia conveniently, thus the methods of embedding this structured information can be easily adapted to various domains.

The major contributions of this work are concluded to three folds: 1) We formulate a novel Hierarchical Semantic Embedding (HSE) framework that integrates semantically structured information of category hierarchy into the deep neural network for FGIR. To our knowledge, this is the first work that explicitly incorporates this structured information to aid FGIR. 2) We introduce a four-level category hierarchy for the Caltech-UCSD birds dataset [39] and construct a new large-scale butterfly dataset that also covers four-level categories for evaluation. To our knowledge, these two datasets are the first that involves in four-level categories in FGIR and they may benefit research on multi-granularity image recognition. 3) We conduct experiments on the two and the VegFru [14] datasets, and demonstrate the effectiveness of our proposed HSE framework over the baseline and existing state-of-the-art methods. Moreover, we also conduct ablative studies to carefully evaluate and analyze the contribution of each component of the proposed framework. *The code, trained models, and dataset are available online: <https://github.com/HCLab-SYSU/HSE>.*

2 RELATED WORK

2.1 Fine-grained image recognition

Recent progress on image classification mainly benefited from the advancement of deep Convolutional Neural Networks (CNNs) [3, 4, 11, 22, 23, 35] that learned powerful feature representation via stacking multiple nonlinear transformations. To adapt the deep CNNs for handling the FGIR task, a bilinear model [25] was proposed to compute high-order image representation that captured local pairwise interactions between features generated by two independent sub-networks, but the bilinear feature is extremely high-dimensional, making it impractical for subsequent analysis. To reduce the feature dimension while keeping comparable performance on FGIR task, Gao et al. [9] developed a compact model that approximates bilinear feature with the polynomial kernels. Kong et al. [19] proposed classifier co-decomposition to further compress the bilinear model.

To better capture subtle visual difference among sub-ordinate categories, a series works [16, 46, 47] were also proposed to leverage extra supervision of bounding boxes and parts to locate discriminative regions. However, the heavy involvement of manual annotations prevents these methods from application to large-scale real-world problems. Recently, visual attention models [5, 26, 30, 42] were intensively proposed to automatically search the informative regions and various works successfully applied this technique to FGIR [8, 17, 28, 50]. Liu et al. [28] formulated a reinforcement learning framework to adaptively glimpse local regions regarding discriminative object parts and trained the framework using a greedy reward strategy with image-level labels. Zheng et al. [50] introduced a multi-attention convolutional neural network that

learned channel grouping for parts localization, and aggregated features from the located regions as well as the global object for classification. These works learned to locate informative regions merely based on image content by the self-attention mechanism. In contrast, some works also introduced extra guidance to learn more meaningful and semantic-related regions to aid FGIR. For example, Liu et al., [2, 27] introduce part-based attribute to guide learning more discriminative features for fine-grained bird recognition. Similarly, He et al. [12] further utilized more detailed language descriptions to help mine discriminative parts or characteristics.

Our framework is also related to some existing works that exploit category hierarchy. For example, Srivastava et al. [37] exploited class hierarchy prior to transfer knowledge among similar lower-level classes for transfer learning. Jia et al. [6] proposed a probabilistic classification model based on a hierarchy and exclusion graph to capture label relations of mutual exclusion, overlap, and subsumption for object classification. Works [5, 41] utilized an RNN to model label co-occurrence dependencies for multi-label recognition. In contrast to these methods that merely model dependencies on label space, our HSE framework introduces the hierarchical information to progressively regularize label prediction and simultaneously guide learning finer-grained feature representation. Besides, using predicted results of the higher level as soft targets for label regularization can distill knowledge learned from the high level to lower level, which is also original compared with these methods.

2.2 Fine-grained image datasets

In the past decade, datasets of FGIR have intensively emerged across various domains ranging from man-made objects to natural plants or animals, including FGVC-Aircraft [29], Stanford Cars [21], Caltech-UCSD birds [39], Stanford Dogs [18], Oxford Flowers [31], to name a few. As a representative dataset that was widely used in previous FGIR works [9, 12, 27], Caltech-UCSD birds dataset contained 11,788 images and covered 200 species of birds. These datasets significantly evolved the research of FGIR, but they primarily focus on categories of one certain level, e.g., Caltech-UCSD birds with 200 species of birds and Stanford Dogs with 120 breeds of dogs. More recently, there also released some datasets that involved categories of multiple levels, like CompCars [45], Boxcars [36], Cars-333 [44] with three-level car categories of make, model, and year, and VegFru [14] with 25 upper-level categories and 292 subordinate classes of vegetables and fruits. These datasets mainly include man-made vehicles [36, 44, 45] and domestic food materials [14]. To better evaluate our proposed frameworks and increase the diversity of dataset with categories of multiple levels, we further organize the 200 bird species with four-level category hierarchy and construct a new butterfly dataset that also covers four-level categories. Besides the research on FGIR with categories of multiple levels, these two datasets have potential to benefit practical applications of wildlife recognition, protection, and discovery.

3 HSE FRAMEWORK

In this section, we describe the proposed HSE framework in detail. Given an image, the framework first utilizes a trunk network to extract image feature maps $f_I \in \mathcal{R}^{W' \times H' \times C'}$, where W' , H' and C' denote the width, height and channel number of the feature maps,

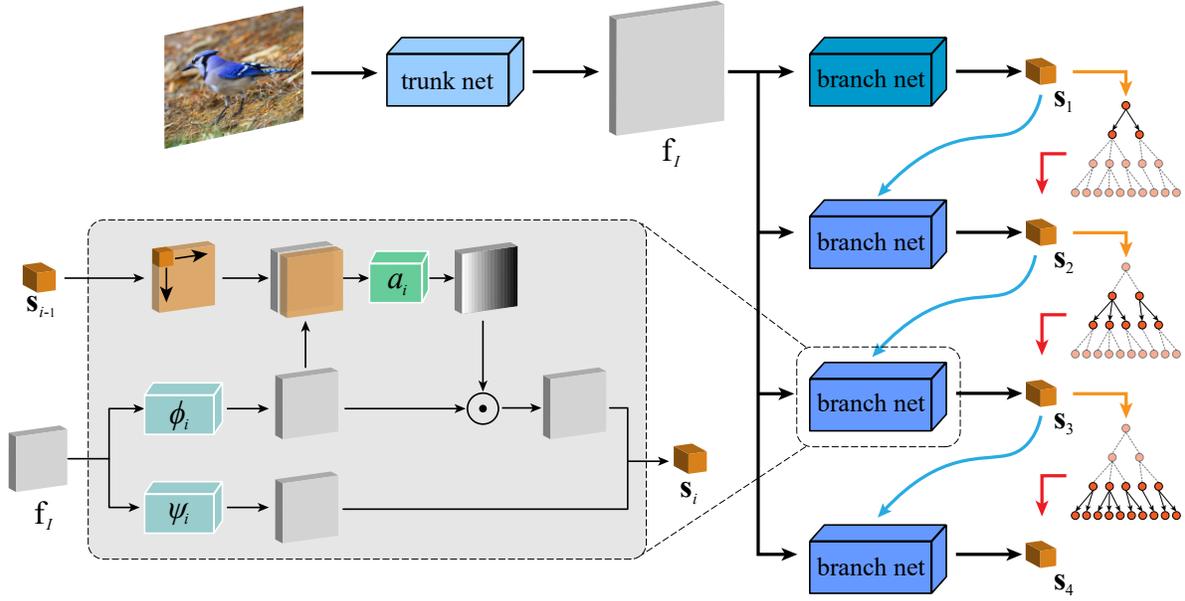


Figure 1: An overall pipeline of our proposed hierarchical semantic embedding framework. It employs a trunk network to extract image features and subsequently utilizes a branch network to predict the categories of each level. At each level, it incorporates the predicted score vector to guide learning finer-grained feature and simultaneously regularizes label prediction during training.

respectively. Then, it orderly utilizes a small branch network to predict the score vectors of all levels, from highest to lowest. At each level, the branch network incorporates the predicted score vector of higher level as prior guidance to learn finer-grained representation via a soft attention mechanism and aggregates this representation with features learned without guidance to predict the score vector of this level. During training, we further use the predicted score vector of higher level as soft targets to regularize the label prediction, such that the predicted result at this level tends to accord with that predicted at the higher level. Since there is no guidance at the first level, we merely use the representation learned without guidance to make prediction and no label regularization is involved either. Fig. 1 gives an overall illustration of the HSE framework.

Before delving deep into the formulation, we first present some notations associated with our task that will be used throughout this article. Without loss of generality, we consider the FGIR task with a category hierarchy of L levels. We utilize l_1, l_2, \dots, l_L to denote each level and s_1, s_2, \dots, s_L to denote the predicted score vectors correspondingly. n_1, n_2, \dots, n_L are used to represent the category number for each level, respectively.

3.1 Semantic embedding representation learning

As we orderly predict the score vector of each level, s_{i-1} is given when making prediction at level l_i . Generally, s_{i-1} encodes the category that the object of the given image belongs to with a high probability at level l_{i-1} , and make prediction at level l_i may tend to distinguish the sub-ordinate categories of this category. As discussed above, some certain parts play key roles to distinguish the

sub-ordinate categories of a superclass. In this work, we take full advantage of this information by incorporating s_{i-1} to guide learning finer-grained feature representation at level l_i . Naturally, this can be implemented by a soft mechanism that learns to fixate on the discriminative regions under the guidance of s_{i-1} .

At level l_i , we first map the image feature maps f_I to higher-level features $\hat{f}_i \in \mathcal{R}^{W \times H \times C}$ via

$$\hat{f}_i = \phi_i(f_I), \quad (1)$$

where $\phi_i(\cdot)$ is a transformation that is implemented by a small network. Then, at each location (w, h) , we introduce a shared attentional mechanism $a_i(\cdot)$ to compute the attention coefficient vector under the guidance of s_{i-1} by

$$\hat{e}_{iwh} = a_i([\hat{f}_{iwh}, \varphi_i(s_{i-1})]), \quad (2)$$

where $\hat{e}_{iwh} = \{\hat{e}_{iwh1}, \hat{e}_{iwh2}, \dots, \hat{e}_{iwhC}\}$ denote the importance of each neuron of feature vector \hat{f}_{iwh} . In the equation, $\varphi_i(\cdot)$ is a linear transformation that transforms s_{i-1} to a semantic feature vector. To make the coefficients easily comparable across different channels, we normalize the coefficients across all the locations of each channels c using a softmax function

$$e_{iwhc} = \frac{\exp(\hat{e}_{iwhc})}{\sum_{w', h'} \exp(\hat{e}_{iwh'c})}. \quad (3)$$

In this way, we can obtain $e_{iwh} = \{e_{iwh1}, e_{iwh2}, \dots, e_{iwhC}\}$ denoting the normalized weight of each neuron of feature vector \hat{f}_{iwh} . Finally, we perform weighted average across all locations of each

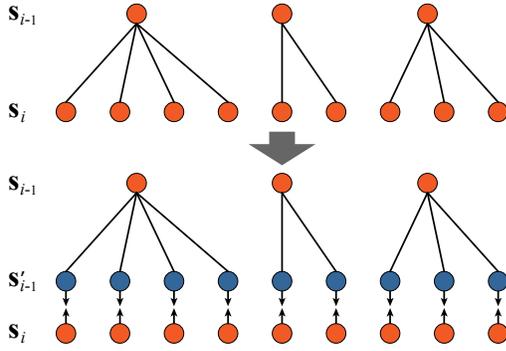


Figure 2: An illustration of the semantic guided label regularization. Top: correlations among categories of level l_{i-1} and l_i . Bottom: s_{i-1} is first extended to s'_{i-1} according to the structured correlations and s_i is pulled close to s'_{i-1} for regularization.

channel to produce the final finer-grained features

$$\mathbf{f}_i = \sum_{w,h} \mathbf{e}_{iwh} \odot \hat{\mathbf{f}}_{iwh}, \quad (4)$$

where \odot denotes the element-wise multiplication operation.

As the feature vector \mathbf{f}_i pays much attention to the local discriminative regions that may tend to capture subtle difference for distinguishing sub-ordinate categories of a superclass. It may ignore the overall description of the object and some background information that may provide contextual cues. Thus, we further extract a feature vector directly from the image feature maps \mathbf{f}_I without guidance for complementary. Similarly, we also adopt a simple transformation $\psi_i(\cdot)$ on \mathbf{f}_I by

$$\hat{\mathbf{f}}'_i = \psi_i(\mathbf{f}_I), \quad (5)$$

where $\hat{\mathbf{f}}'_i \in \mathcal{R}^{W \times H \times C}$. Similar to [11], we simply perform average pooling to obtain the feature vector

$$\mathbf{f}'_i = \frac{1}{WH} \sum_{w,h} \hat{\mathbf{f}}'_{iwh}. \quad (6)$$

The obtained feature vectors \mathbf{f}'_i , \mathbf{f}_i and the concatenation of them $[\mathbf{f}_i, \mathbf{f}'_i]$ are fed to three classifiers to predict the score vectors independently, which are then averaged to produce the final score vector \mathbf{s}_i .

Network details. Similar to recent FGIR works [27, 28], we implement our framework based on the widely used ResNet-50 [11]. Specifically, we implement the trunk network with the preceding 41 convolutional layers of the ResNet-50, and the transformations of $\phi_i(\cdot)$, $\psi_i(\cdot)$ with the following 9 layers of the ResNet-50. We make the trunk network be shared across different levels to better balance prediction accuracy and computational efficiency. $\phi_i(\cdot)$ is simply implemented by a single fully connected layer that map the c -dim score vector to a 1,024-dimension features and the attention mechanism $a_i(\cdot)$ is implemented by two stacked fully connected layers, in which the first one is $c+1,024$ to 1,024 followed by the tanh non-linear function and the second one is 1,024 to c . As we use the identical architecture with ResNet-50, c is 2,048 in this paper.

3.2 Semantic guided label regularization

The hierarchy encodes rich semantic correlations among categories across different levels. For example, the ground truth category at level l_i is the child sub-category of the ground truth category at level l_{i-1} . This correlation information can effectively regularize semantic space and thus make prediction less ambiguous. These correlations should also be maintained among predicted categories of different levels. To this, we incorporate s_{i-1} as soft targets to regularize label prediction at level l_i .

Given the predicted score vector $\mathbf{s}_{i-1} = \{s_{i-1,1}, s_{i-1,2}, \dots, s_{i-1,n_{i-1}}\}$, a high value $s_{i-1,c}$ denotes high confidence that the object in given image belongs to category c at level l_{i-1} , and the predicted scores for the corresponding child sub-categories at level l_i should also be assigned with high values. To this, we first extend \mathbf{s}_{i-1} to \mathbf{s}'_{i-1} according to the structured correlations thus that \mathbf{s}'_{i-1} has the same dimension as \mathbf{s}_i and pull \mathbf{s}_i close to \mathbf{s}'_{i-1} , as shown in Fig. 2. Concretely, if category c at level l_{i-1} has k child sub-categories at level l_i , we duplicate the score $s_{i-1,c}$ by k times. Then we orderly get these duplicated scores together and re-arrange their subscripts to obtain the extended score vector $\mathbf{s}'_{i-1} = \{s'_{i-1,1}, s'_{i-1,2}, \dots, s'_{i-1,n_i}\}$. To make these two vectors easily comparable, we normalize them into probability distribution using the softmax function with temperature T

$$p'_{i-1,c} = \frac{\exp(\frac{s'_{i-1,c}}{T})}{\sum_{c'} \exp(\frac{s'_{i-1,c'}}{T})}, p_{i,c} = \frac{\exp(\frac{s_{i,c}}{T})}{\sum_{c'} \exp(\frac{s_{i,c'}}{T})}, \quad (7)$$

where T is normally set to 1, and we use a high temperature to produce softer probability distribution over classes in our experiment. In this way, we can obtain two normalized probability distributions, i.e., $\mathbf{p}'_{i-1} = \{p'_{i-1,1}, p'_{i-1,2}, \dots, p'_{i-1,n_i}\}$ and $\mathbf{p}_i = \{p_{i,1}, p_{i,2}, \dots, p_{i,n_i}\}$, and define a regularization term as the Kullback-Leibler divergence from \mathbf{p}'_{i-1} to \mathbf{p}_i

$$\ell_i^r = D_{KL}(\mathbf{p}'_{i-1} \parallel \mathbf{p}_i) = - \sum_c p'_{i-1,c} \log \frac{p_{i,c}}{p'_{i-1,c}}. \quad (8)$$

As ℓ_i^r is defined on a single sample, we simply sum up ℓ_i^r over the training set to define the regularization loss term \mathcal{L}_i^r . As suggested in [13], when using soft targets that have high entropy, more information can be provided than hard target per training sample, and the gradient between training samples enjoy less variance. Thus, it can be trained more steadily and using much less training samples. In our experiments, T is set as 4 to produce a sufficiently soft target.

3.3 Optimization

Besides the regularization term, we also employ the cross-entropy loss with the correct labels as the objective function. We first normalize the predicted score vector using exactly the same logits in softmax function but at a normal temperature of 1, expressed as

$$p_{i,c} = \frac{\exp(s_{i,c})}{\sum_{c'} \exp(s_{i,c'})}. \quad (9)$$

Then suppose the ground truth label at level l_i is c_i , its loss can be defined as

$$\ell_i^c = - \sum_c \mathbf{1}(c = c_i) \log p_{i,c}, \quad (10)$$

where $\mathbf{1}(\cdot)$ is the indication function that is assigned as 1 if the expression is true, and assigned as 0 otherwise. We have define the same loss for the score vectors predicted by the three classifier, respectively. Thus, each sample has four losses, and we sum up the four losses over the training set to define the classification loss \mathcal{L}_i^c .

The proposed framework consists of a trunk network and L branch network, and it is trained using a weighted combination of the classification and regularization losses. The training process is empirically divided into two stages, i.e. level-wise training followed by joint fine tuning.

Stage 1: Level-wise training. When training the branch network of level l_i , it needs the predicted score vector of level l_{i-1} to define the regularization loss. Thus, we first train the branch networks in a level-wise manner, from level l_1 to l_L . As our framework is implemented based on the ResNet-50 [11], we initialize the parameters with those of the corresponding layers of ResNet-50 pre-trained on the ImageNet dataset [7]. Concretely, the parameters of the trunk network are initialized by those of the corresponding 41 convolutional layers and the parameters of the transformation $\phi_i(\cdot)$ and $\psi_i(\cdot)$ are initialized with those of the 9 corresponding layers. The parameters of other modules, including the attentional mechanism $a_i(\cdot)$, semantic mapper $\varphi_i(\cdot)$ and the three classifiers, are automatically initialized with the Xavier algorithm [10]. As the trunk network is shared by all branch networks, its parameters are kept fixed at this stage. We train the branch network of level l_i with a weighted combination of the classification and regularization losses

$$\mathcal{L}_i = \mathcal{L}_i^c + \gamma \mathcal{L}_i^r, \quad (11)$$

where γ is a balance parameter. As discussed in [13], the magnitudes of the gradients produced by \mathcal{L}_i^r are scaled by $\frac{1}{T^2}$, thus it is important to multiply them by a scale of T^2 . Thus, we set γ as T^2 , i.e., 16 in our experiments. Note that we merely use the classification loss \mathcal{L}_1^c to train the branch network of level l_1 , as there is no guidance to define the regularization loss term at this level. Similar to previous works [25, 28] on FGIR task, we resize the input images to 512×512 and perform randomly cropping with a size of 448×448 and their horizontal reflections for data augmentation. Then, we train the branch network using the stochastic gradient descent (SGD) algorithm with a batch size of 8, a momentum of 0.9 and a weight decay of 0.00005. The initial learning rate is set as 0.001, and it is divided by 10 when the error plateaus.

Stage 2: Joint fine tuning. After all branch networks are trained, we jointly fine tune the entire framework by combining the loss terms over all granularities

$$\mathcal{L} = \mathcal{L}_1^c + \sum_{i=2}^L \mathcal{L}_i. \quad (12)$$

We adopt the same strategies for data augmentation and hyperparameter setting as Stage 1 except using a smaller initial learning rate 0.0001.

4 DATASETS

We construct a new large-scale butterfly (Butterfly-200) dataset with four-level categories and organize the 200 bird species of the Caltech-UCSD Birds (CUB) dataset also with four-level categories. We evaluate our proposed framework, the baseline methods and

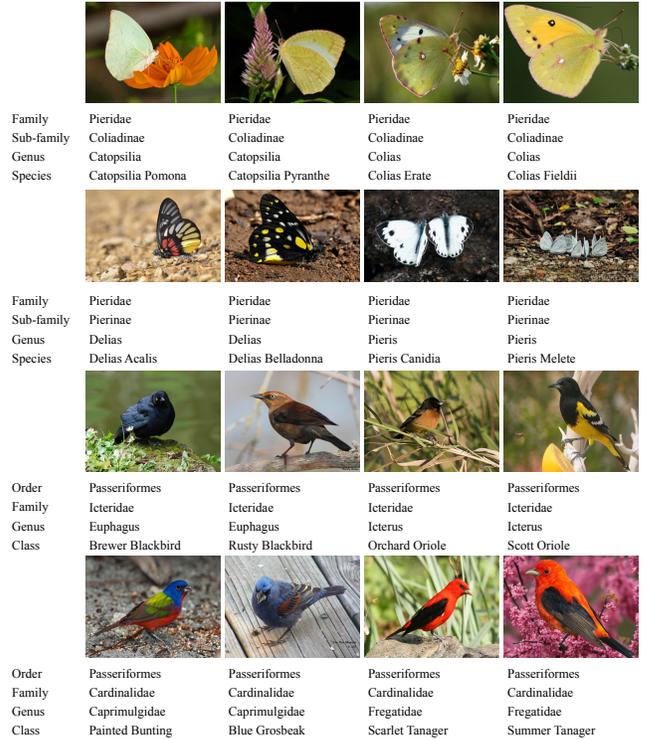


Figure 3: Some samples and their corresponding hierarchical labels from the family of "Pieridae" in the Butterfly-200 dataset (the first two rows) and from the order of "Passeriformes" in the CUB dataset (the last two rows).

the existing competitors on these two and the VegFru [14] datasets. In this section, we first introduce these three datasets.

4.1 Butterfly-200 dataset construction

We select 200 common species of butterflies and build the hierarchical structure with 116 genera, 23 subfamilies, and 5 families according to the insect taxonomy. The butterfly images are collected from two scenarios, natural images with the butterfly in their natural living environment and standard images with the butterfly in the form of specimens, as both are widely used in the real-world applications. The natural images are collected by searching the keywords of butterfly species names on the internet including Google, Flickr, Bing, Baidu, etc. The standard images are collected by capturing the samples in Lab. In this way, a large number of candidate images for each species are collected. To ensure the dataset highly reliable, the candidate images are carefully identified by four experts on butterflies. Currently, we have collected 25,279 butterfly images of the 200 species, with each species containing 30 images at least, which are divided into training, validation, and test set for evaluation. For each species, we randomly select 20% for training, 20% for validation and the rest 60% for test, resulting in a training of 5,135 images, a validation set of 5,135 images, and a test set of 15,009 images, respectively. Figure 3 shows some samples from the family of "Pieridae" and their corresponding hierarchical labels.

Methods	CUB				Butterfly-200			
	l_1 : order	l_2 : family	l_3 : genus	l_4 : species	l_1 : family	l_2 : sub-family	l_3 : genus	l_4 : species
Baseline	98.8	95.0	91.5	85.2	98.9	97.6	94.8	85.1
Baseline+backtrack	98.6	95.1	90.9	85.2	98.7	97.2	94.1	85.1
Ours w/o SERL	98.8	95.1	91.9	86.6	98.9	97.4	95.3	85.8
Ours w/o SGLR	98.8	95.6	92.2	86.7	98.9	97.6	95.1	85.5
Ours (full)	98.8	95.7	92.7	88.1	98.9	97.7	95.4	86.1

Table 1: Comparison of the accuracy (in %) of all levels of our HSE framework, two baseline methods, and two variants of our framework that removes semantic embedding representation learning (Ours w/o SERL) and that removes semantic guided label regularization (Ours w/o SGLR) on the CUB and Butterfly-200 test sets, respectively.

4.2 Caltech-UCSD birds dataset extension

The CUB dataset [39] is the most widely used benchmark for FGIR task. It covers 200 species of birds and contains 11,788 bird images that are divided into a training set of 5,994 images and a test set of 5,794 images. In this work, we build a bird taxonomy hierarchy according to the ornithological systematics, which groups the 200 species into 122 genera, 37 families, and 13 orders. We follow the standard train/test split as [39] for evaluation. Figure 3 also shows some samples from the order of "Passeriformes" and their corresponding hierarchical labels.

4.3 VegFru dataset introduction

VegFru [14] is a newly released large-scale dataset for fine-grained vegetables and fruits recognition. It covers two-level categories of 25 upper-level categories and 292 subordinate classes. The dataset contains 160,731 images in total, including a training set of 29,200 images, a validation set of 14,600 images, and a test set of 116,931 images. Similarly, we follow this standard train/val/test splits as [14] to evaluate our HSE framework and the existing methods for fair comparison.

5 EXPERIMENT

5.1 Significance of semantic embedding

We first implement two baseline methods that use network architecture similar to ours but do not consider the structured correlations to demonstrate the effectiveness of the proposed HSE framework. **Baseline.** Similar to our framework, we utilize a trunk network to extract image features and then utilize four small networks to predict the category of all levels, separately. For fair comparison, we also implement the trunk network with the preceding 41 convolutional layers of the ResNet-50 and the small network with the following 9 layers.

Baseline+backtrack. We utilize the baseline methods to predict the category of the finest level, and backtrack through the hierarchy to obtain the categories of the other levels.

We compare the HSE with these two baseline methods on the CUB and Butterfly-200 datasets in Table 1. Here, we present the accuracies of all levels for comprehensive comparisons. At level l_1 , we find the HSE achieves comparable accuracies with those of the two baseline methods, as there is no semantic guidance at this level. However, at level l_2 to l_4 , the HSE performs consistently better than the baseline methods on both datasets. For example

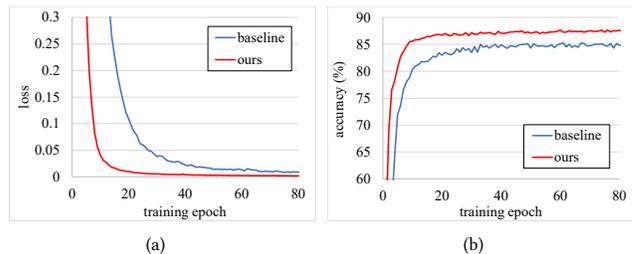


Figure 4: Analysis of the effect of semantic embedding on network learning. These experiments are conducted on categories at level l_4 on the CUB dataset. (a) and (b) are the curves of loss v.s. training epoch on the training set and accuracy v.s. training epoch on the test set, respectively.

on the CUB dataset, the HSE achieves accuracies of 95.7%, 92.7%, and 88.1%, outperforming the baseline methods by 0.6%, 1.2%, and 2.9%, respectively. It is noteworthy that the improvement is more obvious for predicting categories of finer levels, e.g., 1.2% accuracy improvement at level l_3 while 2.9% at level l_4 on the CUB dataset. This phenomenon suggests that incorporating semantic correction information benefits more to challenging tasks.

To delve deep into the effect of semantic embedding on network learning, we further present the curve of loss v.s. training epoch on the training set and the curve of accuracy v.s. training epoch on the test set in Fig 4. These experiments are conducted on recognizing the category of l_4 on the CUB dataset. Compared with the baseline, the HSE can be trained more stably and converged faster.

The foregoing comparisons with the baseline methods demonstrate the effectiveness of the HSE as a whole. Actually, the HSE incorporates the semantic correlation information from two aspects, i.e., semantic embedding representation learning (SERL) and semantic guided label regularization (SGLR). Here, we further conduct ablative studies to assess the actual contributions of these two components.

Contribution of semantic guided label regularization (SGLR). We first evaluate the contribution of SGLR by comparing the performance with and without regularization loss. Specifically, we simply remove the regularization loss terms of each level with others keep fixed and re-train the model in an identical way. As shown in Table 1, removing this term leads to an obvious drop in performance over all levels on both datasets.



Figure 5: Sample number of inter-superclass and intra-superclass errors of our framework with and without SGLR on the (a) CUB and (b) Butterfly-200 datasets.

We further analyze how SGLR improves the performance. When the category of an image is wrongly predicted, we denote it as an inter-superclass error if the wrongly predicted category and ground truth category do not belong to the same superclass, and denote it as an intra-superclass error if they belong to the same superclass. As discussed before, SGLR regularizes label prediction thus that the predicted category at level l_i tends to be the child sub-category of the predicted category at level l_{i-1} . Thus, this tends to help correct the inter-superclass error. To validate this, we present the sample number of inter-superclass and the intra-superclass errors at level l_4 of our HSE with and without SGLR on both datasets. As shown in Fig. 5, introducing SGLR mainly reduces the sample number of inter-superclass error (17.5% relative reduction on the CUB dataset and 13.5% on the Butterfly-200 dataset), finely in accordance with our motivation.

Contribution of semantic embedding representation learning (SERL). Here, we evaluate the benefit of SERL. To this, we remove the feature embedding module (i.e., ϕ_i and a_i) and simply use the feature without guidance for recognition. To ensure fair comparisons, we also re-train the model with both of the classification and regularization losses. Similarly, the performance at each level suffers from an evident drop on both datasets.

As discussed before, SERL helps to attend regions that help to distinguish sub-ordinate categories of the predicted superclass of the higher level. Here, we visualize the attentional regions learned by our HSE framework in Fig. 6. At each row, we present some samples of a specific species, and the first two species belong to the same genus while the last two belong to another genus. For the samples from different species of the same genus, our framework actually attends discriminative regions to better distinguish these species. For example, to differentiate the species of “Bohemian Waxwing” and “Cedar Waxwing” that belong to the genus of “Phoebastria”, the HSE pay much attention to the throat and wing tail regions that provide most discriminative information.

5.2 Comparison with state-of-the-art methods

In this subsection, we compare the HSE framework with existing state-of-the-art methods on the CUB [39] and VegFru [14] datasets. Here, we evaluate on recognizing the categories of the finest level

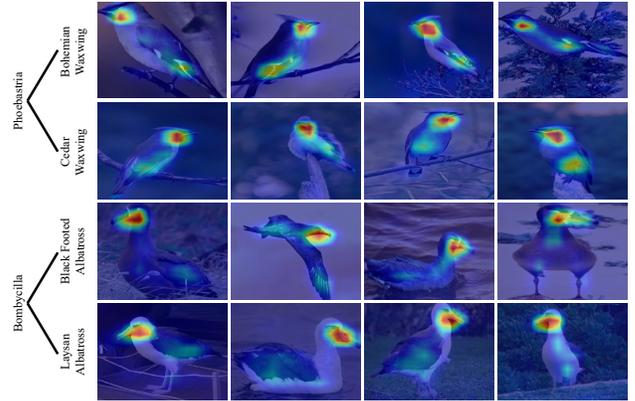


Figure 6: Visualization of the attentional regions learned by the HSE framework. At each row, we present some samples of a specific species, and the first two species belong to the same genus while the last two belong to another genus.

(200 species on CUB and 292 subcategories on VegFru) as existing methods primarily report their results of this level.

Comparison on Caltech-UCSD birds dataset. CUB dataset is the most widely used benchmark for FGIR task, and most works have reported their results on this dataset. We compare our HSE framework with 17 state-of-the-art methods, including Deep Localization, Alignment and Classification (DeepLAC) [24], Semantic Part Detection and Abstraction (SPDA-CNN) [46], Part-RCNN [47], Part Alignment-based (PA-CNN) [20], Pose Normalized CNN (PN-CNN) [1], Picking Deep Filter Responses (PDFR) [48], Multiple Granularity (MG-CNN) [40], Spatial Transformer (ST-CNN) [17], Bilinear-CNN (B-CNN) [25], Compact Bilinear CNN (CB-CNN) [9], Two-Level Attention Network (TLAN) [43], Diverse Attention Network (DAN) [49], Fully Convolutional Attentional Network (FCAN) [28], Recurrent Attention (RA-CNN) [8], Combine Vision and Language (CVL) [12], Attribute-Guided Attention Localization (AGAL) [27], Multi-Attentional CNN (MA-CNN) [50]. Among these methods, some use merely image-level labels (i.e., image-level setting), and some also use bounding box/parts annotations (i.e., box-level setting); thus we also present these information for fair and direct comparisons.

Under the box-level setting, the previous well-performing methods include PN-CNN and B-CNN that achieve accuracies of 85.4% and 85.1%. However, PN-CNN requires strong supervision of both human-defined bounding box and ground truth parts while B-CNN relies on a very high-dimension feature representation (250k dimensions). Under the image-level setting, most works resort to attentional model that automatically search the discriminative regions and aggregate deep features of these regions for classification. For example, MA-CNN learns to attend multiple discriminative regions, and adopt a CNN to extract the global feature from the whole and multiple part-CNNs to extract the local feature from each attentional regions. It achieves an accuracy of 86.5%, which is the best among existing methods. Different from these methods, our HSE framework requires no bounding box and part annotations and does not use multiple CNN to extract local and global features.

Methods	BA	PA	Acc. (%)
Part-RCNN [47]	✓	✓	76.4
DeepLAC [24]	✓	✓	80.3
SPDA-CNN [46]	✓	✓	85.1
PN-CNN [1]	✓	✓	85.4
Part Alignment-CNN [20]	✓		82.8
CB-CNN w/ bbox [9]	✓		84.6
FCAN w/ bbox [28]	✓		84.7
B-CNN w/ bbox [25]	✓		85.1
AGAL w/ bbox [27]	✓		85.5
<hr/>			
TLAN [43]			77.9
DVAN [49]			79.0
MG-CNN [40]			81.7
B-CNN w/o bbox [25]			84.1
ST-CNN [17]			84.1
FCAN w/o bbox [28]			84.3
PDFR [48]			84.5
CB-CNN w/o bbox [9]			85.0
RA-CNN [8]			85.3
AGAL w/o bbox [27]			85.4
CVL [12]			85.6
MA-CNN [50]			86.5
<hr/>			
Ours			88.1

Table 2: Comparisons of our HSE framework with existing state of the arts on recognizing categories of level l_4 on the CUB dataset. BA and PA denote bounding box annotations and part annotations, respectively. ✓ indicates corresponding annotations are used during training or test.

Instead, it embeds structure information of category hierarchy to learn fine-grained feature representation and regularize label prediction, leading to obvious performance improvement, i.e., 88.1% in accuracy.

Note that our HSE introduces extra guidance of the category hierarchy. However, this hierarchy can be easily obtained from the literature of taxonomy or retrieved from the Wikipedia. Besides, we also compare with existing methods that also rely on extra supervisions, like AGAL requiring attribute annotations and CVL depending on sentence description. Our HSE achieves an accuracy of 88.1%, much better than theirs, i.e., 85.5% and 85.6%, respectively. **Comparison on VegFru dataset.** VegFru is a newly released large-scale dataset for fine-grained vegetables and fruits recognition, and some works also report their results on this dataset. Here, we also present comparisons with the baseline and existing methods on this dataset in Table 3. As shown, the HSE also significantly outperforms all these methods.

6 CONCLUSION

Fine-grained categories naturally form a hierarchy with different levels of concept abstraction, and this hierarchy encodes rich correlations among categories across different levels. In this work, we investigate simultaneously predicting categories of all levels in the hierarchy and integrating this structured correlation information into the deep neural network by developing a novel Hierarchical

Methods	Acc. (%)
Baseline	87.1
CB-CNN [9]	82.2
HybridNet [14]	83.5
Ours (full)	89.4

Table 3: Comparison of accuracy of our HSE framework, existing state-of-the-art methods, and the baseline methods on the VegFru dataset.

Semantic Embedding (HSE) framework. Specifically, the HSE orderly predicts the score vector for each level, and at each level, it incorporates the predicted score vector of the higher level to guide learning finer-grained feature representation and simultaneously regularize label prediction during training. To evaluate the HSE framework, we extend the Caltech-UCSD birds with four-level categories and construct a butterfly dataset also with four-level categories. Extensive experiments and thorough analysis on these two and the VegFru datasets demonstrate the superiority of the proposed HSE framework over the baseline methods and existing competitors.

ACKNOWLEDGEMENT

We would like to thank Prof. Min Wang, Associate Prof. XiaoLing Fan, Dr. Haiming Xu, and Dr. Hailing Zhuang with Department of Entomology, College of Agriculture, South China Agricultural University for their assistance in butterfly image annotations. This work was supported in part by the Chinese National Science Foundation (NSFC No. 61702196), and Science and Technology Planning Project of Guangdong Province, China (No. 2017A020208041). This work is jointly supported by State Key Development Program under Grant 2018YFC0830103.

REFERENCES

- [1] Steve Branson, Grant Van Horn, Serge Belongie, and Pietro Perona. 2014. Bird species categorization using pose normalized deep convolutional nets. *British Machine Vision Conference (2014)*.
- [2] Tianshui Chen, Liang Lin, Riquan Chen, Yang Wu, and Xiaonan Luo. 2018. Knowledge-Embedded Representation Learning for Fine-Grained Image Recognition. In *Proc. of International Joint Conference on Artificial Intelligence*. 627–634.
- [3] Tianshui Chen, Liang Lin, Lingbo Liu, Xiaonan Luo, and Xuelong Li. 2016. DISC: Deep Image Saliency Computing via Progressive Representation Learning. *IEEE Trans. Neural Netw. Learning Syst.* 27, 6 (2016), 1135–1149.
- [4] Tianshui Chen, Liang Lin, Wangmeng Zuo, Xiaonan Luo, and Lei Zhang. 2018. Learning a Wavelet-like Auto-Encoder to Accelerate Deep Neural Networks. In *Proc. of AAAI Conference on Artificial Intelligence*. 6722–6729.
- [5] Tianshui Chen, Zhouxia Wang, Guanbin Li, and Liang Lin. 2018. Recurrent Attentional Reinforcement Learning for Multi-label Image Recognition. In *Proc. of AAAI Conference on Artificial Intelligence*. 6730–6737.
- [6] Jia Deng, Nan Ding, Yangqing Jia, Andrea Frome, Kevin Murphy, Samy Bengio, Yuan Li, Hartmut Neven, and Hartwig Adam. 2014. Large-scale object classification using label relation graphs. In *European conference on computer vision*. Springer, 48–64.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 248–255.
- [8] Jianlong Fu, Heliang Zheng, and Tao Mei. 2017. Look closer to see better: recurrent attention convolutional neural network for fine-grained image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [9] Yang Gao, Oscar Beijbom, Ning Zhang, and Trevor Darrell. 2016. Compact bilinear pooling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 317–326.

- [10] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 249–256.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [12] Xiangteng He and Yuxin Peng. 2017. Fine-grained Image Classification via Combining Vision and Language. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognitions* (2017).
- [13] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).
- [14] Saihui Hou, Yushan Feng, and Zilei Wang. 2017. VegFru: A Domain-Specific Dataset for Fine-grained Visual Categorization. In *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 541–549.
- [15] Qichang Hu, Huibing Wang, Teng Li, and Chunhua Shen. 2017. Deep CNNs With Spatially Weighted Pooling for Fine-Grained Car Recognition. *IEEE Transactions on Intelligent Transportation Systems* 18, 11 (2017), 3147–3156.
- [16] Shaoli Huang, Zhe Xu, Dacheng Tao, and Ya Zhang. 2016. Part-stacked CNN for fine-grained visual categorization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1173–1182.
- [17] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. 2015. Spatial transformer networks. In *Advances in Neural Information Processing Systems*. 2017–2025.
- [18] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. 2011. Novel dataset for fine-grained image categorization: Stanford dogs. In *Proc. CVPR Workshop on Fine-Grained Visual Categorization (FGVC)*, Vol. 2. 1.
- [19] Shu Kong and Charles Fowlkes. 2016. Low-rank Bilinear Pooling for Fine-Grained Classification. *arXiv preprint arXiv:1611.05109* (2016).
- [20] Jonathan Krause, Hailin Jin, Jianchao Yang, and Li Fei-Fei. 2015. Fine-grained recognition without part annotations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5546–5555.
- [21] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 2013. 3d object representations for fine-grained categorization. In *Computer Vision Workshops (ICCVW)*, 2013 *IEEE International Conference on*. IEEE, 554–561.
- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
- [23] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [24] Di Lin, Xiaoyong Shen, Cewu Lu, and Jiaya Jia. 2015. Deep lac: Deep localization, alignment and classification for fine-grained recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1666–1674.
- [25] Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. 2015. Bilinear cnn models for fine-grained visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 1449–1457.
- [26] Lingbo Liu, Hongjun Wang, Guanbin Li, Wanli Ouyang, and Liang Lin. 2018. Crowd Counting using Deep Recurrent Spatial-Aware Network. In *Proc. of International Joint Conference on Artificial Intelligence*.
- [27] Xiao Liu, Jiang Wang, Shilei Wen, Errui Ding, and Yuanqing Lin. 2017. Localizing by Describing: Attribute-Guided Attention Localization for Fine-Grained Recognition. In *AAAI*. 4190–4196.
- [28] Xiao Liu, Tian Xia, Jiang Wang, and Yuanqing Lin. 2016. Fully convolutional attention localization networks: Efficient attention localization for fine-grained recognition. *arXiv preprint arXiv:1603.06765* (2016).
- [29] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. 2013. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151* (2013).
- [30] Volodymyr Mnih, Nicolas Heess, Alex Graves, et al. 2014. Recurrent models of visual attention. In *Advances in neural information processing systems*. 2204–2212.
- [31] Maria-Elena Nilsback and Andrew Zisserman. 2008. Automated flower classification over a large number of classes. In *Computer Vision, Graphics & Image Processing*, 2008. *ICVGIP'08. Sixth Indian Conference on*. IEEE, 722–729.
- [32] JV Remsen Jr, Alexis FLA Powell, Richard Schodde, F Keith Barker, and Scott M Lanyon. 2016. A revised classification of the Icteridae (Aves) based on DNA sequence data. *Zootaxa* 4093, 2 (2016), 285–292.
- [33] Rodrigo B Salvador, Henk Van der Jeugd, and Barbara M Tomotani. 2017. Taxonomy of the European Pied Flycatcher *Ficedula hypoleuca* (Aves: Muscicapidae). *Zootaxa* 4291, 1 (2017), 171–182.
- [34] HEMCHANDRANAUTH SAMBHU and ALLIEA NANKISHORE. 2018. Butterflies (Lepidoptera) of Guyana: A compilation of records. *Zootaxa* 4371, 1 (2018), 1–187.
- [35] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [36] Jakub Sochor, Adam Herout, and Jiri Havel. 2016. Boxcars: 3d boxes as cnn input for improved fine-grained vehicle recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3006–3015.
- [37] Nitish Srivastava and Ruslan R Salakhutdinov. 2013. Discriminative transfer learning with tree-based priors. In *Advances in Neural Information Processing Systems*. 2094–2102.
- [38] Rudi Verovnik and Miloš Popović. 2013. Annotated checklist of Albanian butterflies (Lepidoptera, Papilionoidea and Hesperioidea). *ZooKeys* 323 (2013), 75.
- [39] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. 2011. The caltech-ucsd birds-200-2011 dataset. (2011).
- [40] Dequan Wang, Zhiqiang Shen, Jie Shao, Wei Zhang, Xiangyang Xue, and Zheng Zhang. 2015. Multiple granularity descriptors for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision*. 2399–2406.
- [41] Jiang Wang, Yi Yang, Junhua Mao, Zhiheng Huang, Chang Huang, and Wei Xu. 2016. Cnn-rnn: A unified framework for multi-label image classification. In *Computer Vision and Pattern Recognition (CVPR)*, 2016 *IEEE Conference on*. IEEE, 2285–2294.
- [42] Zhouxia Wang, Tianshui Chen, Guanbin Li, Ruijia Xu, and Liang Lin. 2017. Multi-label Image Recognition by Recurrently Discovering Attentional Regions. In *IEEE International Conference on Computer Vision*. IEEE, 464–472.
- [43] Tianjun Xiao, Yichong Xu, Kuiyuan Yang, Jiaying Zhang, Yuxin Peng, and Zheng Zhang. 2015. The application of two-level attention models in deep convolutional neural network for fine-grained image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 842–850.
- [44] Saining Xie, Tianbao Yang, Xiaoyu Wang, and Yuanqing Lin. 2015. Hyper-class augmented and regularized deep learning for fine-grained image classification. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015).
- [45] Linjie Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. 2015. A large-scale car dataset for fine-grained categorization and verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3973–3981.
- [46] Han Zhang, Tao Xu, Mohamed Elhoseiny, Xiaolei Huang, Shaoting Zhang, Ahmed Elgammal, and Dimitris Metaxas. 2016. Spda-cnn: Unifying semantic part detection and abstraction for fine-grained recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1143–1152.
- [47] Ning Zhang, Jeff Donahue, Ross Girshick, and Trevor Darrell. 2014. Part-based R-CNNs for fine-grained category detection. In *European conference on computer vision*. Springer, 834–849.
- [48] Xiaopeng Zhang, Hongkai Xiong, Wengang Zhou, Weiyao Lin, and Qi Tian. 2016. Picking deep filter responses for fine-grained image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1134–1142.
- [49] Bo Zhao, Xiao Wu, Jiashi Feng, Qiang Peng, and Shuicheng Yan. 2017. Diversified Visual Attention Networks for Fine-Grained Object Classification. *IEEE Transactions on Multimedia* 19, 6 (2017), 1245–1256.
- [50] Heliang Zheng, Jianlong Fu, Tao Mei, and Jiebo Luo. 2017. Learning multi-attention convolutional neural network for fine-grained image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5209–5217.