

Deep Learning for Multilabel Remote Sensing Image Annotation With Dual-Level Semantic Concepts

Panpan Zhu, Yumin Tan, Liqiang Zhang[✉], Yuebin Wang[✉], Jie Mei[✉], Hao Liu, and Mengfan Wu

Abstract—Multilabel remote sensing (RS) image annotation is a challenging and time-consuming task that requires a considerable amount of expert knowledge. Most existing RS image annotation methods are based on handcrafted features and require multistage processes that are not sufficiently efficient and effective. An RS image can be assigned with a single label at the scene level to depict the overall understanding of the scene and with multiple labels at the object level to represent the major components. The multiple labels can be used as supervised information for annotation, whereas the single label can be used as additional information to exploit the scene-level similarity relationships. By exploiting the dual-level semantic concepts, we propose an end-to-end deep learning framework for object-level multilabel annotation of RS images. The proposed framework consists of a shared convolutional neural network for discriminative feature learning, a classification branch for multilabel annotation and an embedding branch for preserving the scene-level similarity relationships. In the classification branch, an attention mechanism is introduced to generate attention-aware features, and skip-layer connections are incorporated to combine information from multiple layers. The philosophy of the embedding branch is that images with the same scene-level semantic concepts should have similar visual representations. The proposed method adopts the binary cross-entropy loss for classification and the triplet loss for image embedding learning. The evaluations on three multilabel RS image data sets demonstrate the effectiveness and superiority of the proposed method in comparison with the state-of-the-art methods.

Index Terms—Attention mechanism, dual-level semantic concepts, remote sensing (RS) image multilabel annotation, triplet loss.

I. INTRODUCTION

AUTOMATIC image annotation has been an active research topic in the remote sensing (RS) field. Image annotation is to assign one or several predefined semantic

Manuscript received April 19, 2019; revised August 24, 2019, October 27, 2019, and November 22, 2019; accepted December 9, 2019. This work was supported by the National Natural Science Foundation of China under Grant 41925006 and Grant 41801241. (*Corresponding authors:* Yumin Tan; Liqiang Zhang.)

Panpan Zhu, Liqiang Zhang, Hao Liu, and Mengfan Wu are with the State Key Laboratory of Remote Sensing Science, Faculty of Geographical Science, Beijing Normal University, Beijing 100875, China (e-mail: ppzhu@mail.bnu.edu.cn; zhanglq@bnu.edu.cn; apprentice_g@163.com; wmf1991yeah@126.com).

Yumin Tan is with the School of Transportation Science and Engineering, Beihang University, Beijing 100191, China (e-mail: tanym@buaa.edu.cn).

Yuebin Wang is with the School of Land Science and Technology, China University of Geosciences, Beijing 100083, China (e-mail: xxgedxwyb@163.com).

Jie Mei is with the College of Computer Science, Nankai University, Tianjin 300071, China (e-mail: meijie0507@gmail.com).

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TGRS.2019.2960466

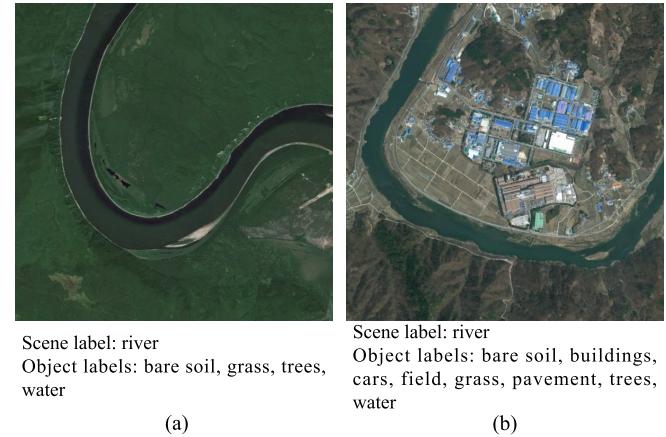


Fig. 1. Illustration of RS images with dual-level semantic concepts. At the scene level, each image is associated with a single label, while at the object level, each image is associated with multiple labels. (a) Scene label: river. Object labels: bare soil, grass, trees, and water. (b) Scene label: river. Object labels: bare soil, buildings, cars, field, grass, pavement, trees, and water.

concepts to an image [1]. Most previous RS image recognition and retrieval methods assume that each image is annotated by a single label, which refers to the overall understanding of the scene or the most significant semantic class in the image [2], [3]. Although such scene labels offer a holistic understanding of the image content, it is insufficient for them to delineate object primitives that exist in the image [4]–[6]. To solve this problem, many semantic segmentation algorithms have been proposed, in which each pixel is assigned with an object label [7], [8]. Semantic segmentation provides precise localization information of objects that exist in an image [9] and has been applied to a wide variety of tasks including road extraction [8] and building footprint extraction [10]. However, the ground truths of pixelwise annotation are mostly obtained by tedious and time-consuming manual labeling [5]. In this case, multilabel annotation has received significant attention because it tells us about objects that exist in the image, while the cost of ground-truth annotations is relatively low [6]. Both scene and object level labels are image level labels (see Fig. 1), which are much easier to obtain compared with pixel level labels required by semantic segmentation.

Multilabel annotation of RS images can be applied to many real applications, such as precise land use/land cover investigation [11] and riverbank buildings monitoring. Take riverbank buildings monitoring, for example, multilabel annotation helps provide important cues of the existence of buildings along the riverside. Riverside buildings are susceptible to rising water

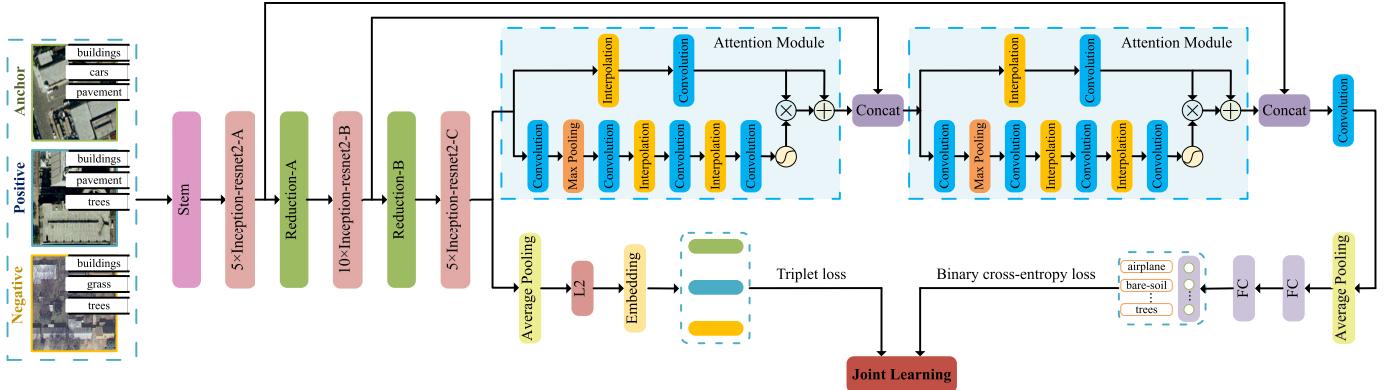


Fig. 2. Proposed RS image annotation framework using dual-level semantic concepts. Three components are included: visual feature learning, multilabel classification, and image embedding learning.

levels. Therefore, obtaining accurate building information is very important not only for riparian environment protection but also for dealing with emergencies such as floods [12]. Through scene classification, Fig. 1(a) and (b) are assigned with scene label *river*. Through multilabel annotation, Fig. 1(a) is assigned with object labels *bare soil*, *grass*, *trees*, and *water*, and Fig. 1(b) is assigned with object labels *bare soil*, *buildings*, *cars*, *field*, *grass*, *pavement*, *trees*, and *water*. The identification of buildings in Fig. 1(b) provides decision-making basis for managers and is conducive to subsequent target detection and instance segmentation.

A. Challenges of Multilabel Annotation

Much effort has been devoted to developing automatic multilabel annotation methods. In [11], classic multilabel learning methods, such as the binary relevance (BR) [13] and multilabel K -nearest neighbor (ML-KNN) [14] methods, were successfully applied to land cover multilabel classification. In [15], a graph-based multilabel annotation method was proposed based on low rank representation. These methods have achieved limited performance due to the utilization of hand-crafted features that cannot express high-level semantics [16]. In contrast, deep learning methods adaptively learn image features so that more discriminative semantic features can be learned [6].

Deep learning methods, in particular, convolutional neural networks (CNNs), have been widely used in the RS field, such as RS image classification [17]–[19] and semantic labeling [20]. For example, Maggiori *et al.* [20] used a specific CNN module to combine features from multiple layers for aerial image segmentation. Gardner and Nichols [21] applied binary-cross entropy loss instead of softmax loss in a CNN model for RS image multilabel annotation. However, the label correlations were ignored in their article. How to utilize label correlations to improve annotation performance is one of the challenges of multilabel annotation [13]. For natural image annotation [22]–[24], the semantic label structure was used to enhance the image annotation performance. The study in [25] and [26] demonstrated that hierarchical label structures could improve the multilabel annotation performance. Therefore, the goal of this article is to explore how to utilize the label structure of RS images to enhance the annotation performance.

B. Contributions of Our Work

In this article, we propose an end-to-end deep learning framework for multilabel annotation of RS images that exploit the dual-level semantic concepts. The framework is shown in Fig. 2. The framework consists of a visual feature learning component, a multilabel classification component and an image embedding learning component. In the visual feature learning component, the Inception-ResNet-v2 architecture is leveraged to extract discriminative visual features. In the multilabel classification component, the semantic concepts at the object level are used as supervised information. Two attention modules are introduced for salient object detection. The skip-layer connections combine fine appearance information from bottom layers and coarse semantic information from top layers. In addition, the binary cross-entropy loss is applied to guide the classification process. In the embedding branch, the similarity relationships of the scene-level semantic concepts are exploited. The triplet loss is adopted for the embedding learning process. The main contributions of our work are summarized as follows.

- 1) We propose an end-to-end deep learning multilabel annotation method. The proposed method takes the dual-level semantic label information of RS images into consideration, which has higher performance compared with the state-of-the-art methods.
- 2) The attention module and skip-layer connections are incorporated to enhance the discriminative abilities of feature representations.
- 3) The triplet loss is leveraged to preserve the similarity relationships of semantic concepts at the scene level.

II. RELATED WORKS

In this section, related works of multilabel image annotation are reviewed. Multilabel image annotation/classification algorithms can be broadly divided into two categories, namely problem transformation methods and algorithm adaptation methods [13]. The underlying principle of problem transformation methods is to convert the original multilabel annotation problem into other well-established learning problems. One of the earliest representative algorithms is the BR approach [13]. In this approach, the multilabel classification problem is

transformed into a multiple independent single label binary classification problem. The inherent drawback of the BR approach is that the relationships among the labels are ignored. ML-KNN [14] and multilabel canonical correlation analysis [27] are typical representatives of label adaptation algorithms. In this article, we follow the idea of problem transformation by transforming a multilabel annotation problem into a multiple binary classification problem, but we additionally exploit the label relationships in an added embedding branch.

A. Exploiting Multilevel Semantic Concepts

Multilevel semantic concepts have been adopted to improve the performance of image annotation in many studies. Hu *et al.* [28] proposed a multilevel max-margin discriminative analysis framework for the annotation of high-resolution RS images. In this article, the multilevel semantics from coarse to fine is taken into consideration. Dumitru *et al.* [29] developed a three-level annotation scheme for high-resolution synthetic aperture radar images. Hu *et al.* [26] proposed a structured model for multilabel image annotation that leverages the hierarchical levels of concepts, including tags, groups, and labels, to encode the semantic correlation. Zhang *et al.* [30] considered four semantic levels, including visual features, object categories, spatial object patterns, and zone functions, to map urban functional zones. In our work, the semantic information at both the scene level and the object level are fully exploited.

B. Attention Mechanism and Skip-Layer Connections

Much effort has been devoted to applying attention mechanisms to deep neural networks [31]–[33]. The underlying principle is to guide the network to focus on the salient objects. For example, the attention mechanism has been encoded into a unified framework of CNN and recurrent neural network (RNN) [22] implicitly to better predict small objects that need more contextual information. In [32], self-attention mechanism was performed to capture long-range dependencies via nonlocal operations. Hu *et al.* [33] introduced a powerful channelwise attention mechanism, termed the Squeeze-and-Excitation block, which explicitly models the interdependencies among channels. In [34], skip-layer connections are used to combine feature maps of different layers. In [35], the attention residual learning structure was used to improve the image classification performance. This structure has two advantages. On the one hand, due to the filtering role of the attention module, the good properties of the original features can be kept. On the other hand, the network has a channel to bypass the soft mask branch to forward features to top layers and weaken the feature selection ability of the mask branch. Borrowing the idea of attention residual learning, we introduce a similar attention mechanism in our proposed method to generate more discriminative feature representations.

C. Image Recognition Using the Triplet Loss

The triplet loss has been used to guide image recognition and clustering, especially face image recognition.

Schroff *et al.* [36] employed a triplet loss function based on the large margin nearest-neighbor method [37] to train a deep CNN for facial image embedding. The idea behind their work is that facial images of the same person usually have smaller distances than those of different persons. Cheng *et al.* [38] proposed a novel multichannel part-based CNN for person re-identification, with an improved triplet loss function guided. In [24], a triplet loss framework was conducted for the learning of a joint embedding space for the images and their tags, where the distances between the images and their tags and that among the images sharing all or part of the common tags were minimized. Inspired by these studies, we include the triplet loss in our framework to guide the learning of image embeddings.

III. PROPOSED FRAMEWORK

In this section, we present our proposed annotation framework in detail. Our framework is designed for object-level multilabel annotation of RS images and consists of three components: a shared CNN for learning the visual representations, a classification branch for the multilabel semantic annotation at the object level and an embedding branch for leveraging the similarity relationships among the images according to the scene-level semantic concepts. Two attention modules are incorporated into the classification branch to generate more discriminative feature representations. Combining the features from bottom layers with features from top layers forms robust features [7]. The entire framework is trained in an end-to-end fashion by jointly optimizing the binary cross-entropy loss and the triplet loss.

A. Problem Formulation

Multilabel annotation aims to assign a set of appropriate labels to unknown samples. Assume that we have an archive $\mathbf{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m\}$ of m images, and $T = \{t^1, t^2, \dots, t^q\}$ denotes the label space with q possible unique labels at the object level. If \mathbf{X}_i contains the j th label, then $t_i^j = 1$; otherwise, $t_i^j = 0$. In addition, each image is associated with a single label at the scene level in the label space $Y = \{y^1, y^2, \dots, y^k\}$. If \mathbf{X}_i contains the j th label, then $y_i^j = 1$; otherwise, $y_i^j = 0$. The scene-level label space is distinguished from the object-level label space in that all the elements in the former space are one-hot representations; in other words, one image may have one or more labels at the object level but only one label at the scene level. The goal of the proposed method is to learn a function, $\mathcal{L} : \mathbf{X} \rightarrow 2^T$, from the training set, $D = \{(\mathbf{X}_i, T_i^g, Y_i) | 1 \leq i \leq m\}$, which includes the images, multiple labels at the object level and single label at the scene level.

B. Visual Feature Learning

We take the Inception-ResNet-v2 [39] model as the shared CNN to extract the visual image features. Since the RS images can be annotated with dual-level semantic concepts, the inception-style CNN architecture [40] with several different filter sizes at one stage is adapted to the dual-level

label characteristics of the RS images. Briefly, the Inception-ResNet-v2 model incorporates the residual connections into the Inception V3 model [40]. We select the output of the last block of the Inception-ResNet-v2-A unit and the Inception-ResNet-v2-B unit to fuse the different-layer features. The feature maps generated from the two blocks have 35×35 and 17×17 pixels, along with 320 and 1088 channels, respectively.

C. Classification Branch

After feature learning, our proposed method is divided into two branches: the classification branch and the embedding branch. The classification branch is composed of two attention modules, a skip-layer connection, a convolution layer, an average pooling layer, two fully connected layers, and a classifier layer. In addition, the binary cross-entropy loss is used for the training of multilabel annotation.

1) *Skip-Layer Connection*: Due to a series of convolution and pooling operations in the feature leaning stage, the feature maps at top layers are coarse and have a low resolution. The 8×8 pixel feature maps at the final Inception-ResNet2-C unit contain coarse semantic information and lose the fine appearance information. To address this, skip-layer connections [7], [41] are added to combine feature maps of bottom layers and top layers. Since the bottom layers retain fine appearance information and the top layers represent the abstract semantic information, skip-layer connections of the bottom and top layers are useful for refining the network [34] and improving the accuracy of the multilabel annotation.

The first attention module upsamples the feature maps into a size of 17×17 pixels with 128 channels. Then, the output is concatenated to the output of the last Inception-ResNet2-B unit by depth to obtain feature maps with a size of $17 \times 17 \times 1216$. Similarly, the second attention module upsamples the feature maps to have a size of 35×35 with 64 channels. After concatenation with the output of the last Inception-ResNet2-A unit, the generated feature maps have a size of $35 \times 35 \times 384$.

2) *Trunk-and-Mask Attention Module*: The trunk-and-mask attention module is inspired by the residual attention network [35]. Briefly, the attention module is composed of two subbranches, i.e., the trunk subbranch and the soft mask subbranch, as shown in Fig. 2. The details of the attention modules are outlined in Table I.

The trunk subbranch includes a bilinear interpolation to upsample the feature maps and a 1×1 convolution layer to reduce dimension. The soft mask subbranch consists of a downsampling operation and two upsampling operations. The downsampling operation is performed by max pooling with a stride of two, and the upsampling is performed by bilinear interpolation. In addition, 1×1 convolutional layers are used for dimension reduction. Finally, a sigmoidal layer is employed to normalize the output range to $(0, 1)$.

Given the input feature map $g(\mathbf{X})$, the output of the trunk subbranch is denoted as $t(g(\mathbf{X}))$, and the output of the soft mask subbranch is denoted as $m(g(\mathbf{X}))$. In this case, the output of the attention module is formulated as

$$h_{i,c} = (1 + m_{i,c}(g(\mathbf{X}))) * t_{i,c}(g(\mathbf{X})) \quad (1)$$

TABLE I
DETAILS OF THE ATTENTION MODULE

	Layer	Output Size	Patch Size / Stride
Attention Module	Trunk Subbranch	Interpolation	$17 \times 17 \times 1, 536$
	Trunk Subbranch	Convolution	$17 \times 17 \times 128$
	Soft Mask Subbranch	Convolution	$8 \times 8 \times 1, 024$
	Soft Mask Subbranch	Max pooling	$4 \times 4 \times 1, 024$
	Soft Mask Subbranch	Convolution	$4 \times 4 \times 128$
	Soft Mask Subbranch	Interpolation	$8 \times 8 \times 128$
	Soft Mask Subbranch	Convolution	$8 \times 8 \times 16$
	Soft Mask Subbranch	Interpolation	$17 \times 17 \times 16$
	Soft Mask Subbranch	Convolution	$17 \times 17 \times 128$
	Concatenation		$17 \times 17 \times 1, 216$
Attention Module	Trunk Subbranch	Interpolation	$35 \times 35 \times 1, 216$
	Trunk Subbranch	Convolution	$35 \times 35 \times 64$
	Soft Mask Subbranch	Convolution	$17 \times 17 \times 512$
	Soft Mask Subbranch	Max pooling	$8 \times 8 \times 512$
	Soft Mask Subbranch	Convolution	$8 \times 8 \times 64$
	Soft Mask Subbranch	Interpolation	$17 \times 17 \times 64$
	Soft Mask Subbranch	Convolution	$17 \times 17 \times 8$
	Soft Mask Subbranch	Interpolation	$35 \times 35 \times 8$
	Soft Mask Subbranch	Convolution	$35 \times 35 \times 64$
	Concatenation		$35 \times 35 \times 384$
	Convolution		$17 \times 17 \times 128$
			$3 \times 3 / 2$

where i ranges over all spatial positions and c ranges over all channels, $m(\mathbf{X}) \in (0, 1)$.

3) *Binary Cross-Entropy Loss*: Image labels provide supervised information for classification. In the classification branch, the binary cross-entropy loss is employed, which is formulated as

$$L_c = \frac{1}{mq} \sum_{i=1}^m \sum_{j=1}^q (t_i^j \log \hat{t}_i^j + (1 - t_i^j) \log (1 - \hat{t}_i^j)) \quad (2)$$

where $t_i^j \in (0, 1)$ denotes the j th ground-truth label for training image \mathbf{X}_i . \hat{t}_i^j is the output of the sigmoidal layer. m is the number of training images and q is the cardinality of the object-level label space.

D. Embedding Branch

Considering that learning embeddings for multilabel image annotation is effective [24], [36], [42], we adopt a triplet loss similar to that used in FaceNet [36]. The embedding branch acts as a constraint to maintain the scene-level semantic similarity relationships among the images. The design of our embedding branch rests on the assumption that images with the same semantic concept at the scene level should have similar visual representations, whereas images with different scene-level semantic concepts should have dissimilar visual representations. Specifically, our embedding branch includes an average pooling layer, an L2 normalization layer, and

an embedding layer. Sections III-D1 and III-D2 describe the triplet loss and how image triplets are selected.

1) *Triplet Loss*: We represent the embedding of the image \mathbf{X} as $f(\mathbf{X}) \in \mathbb{R}^d$, which embeds \mathbf{X} into a d -dimensional Euclidean space. Moreover, the embedding is forced to rest on the d -dimensional hypersphere, i.e., $\|f(\mathbf{X})\|_2 = 1$. An image \mathbf{X}_i^a (anchor) is closer to \mathbf{X}_i^p (positive), which has the same semantic concept, than to \mathbf{X}_i^n , which has a different semantic concept. Thus

$$\begin{aligned} \|f(\mathbf{X}_i^a) - f(\mathbf{X}_i^p)\|_2^2 + \alpha &< \|f(\mathbf{X}_i^a) - f(\mathbf{X}_i^n)\|_2^2 \\ \forall(f(\mathbf{X}_i^a), f(\mathbf{X}_i^p), f(\mathbf{X}_i^n)) \in \mathcal{T} \end{aligned} \quad (3)$$

where α refers to the enforced margin between the positive and negative image pairs. \mathcal{T} denotes the set of all possible triplets in the training set.

The triplet loss is formulated as

$$L_t = \frac{1}{3N} \sum_i^N [\|f(\mathbf{X}_i^a) - f(\mathbf{X}_i^p)\|_2^2 - \|f(\mathbf{X}_i^a) - f(\mathbf{X}_i^n)\|_2^2 + \alpha]_+ \quad (4)$$

where $[x]_+$ represents $\max(x, 0)$ and N refers to the number of triplets.

2) *Triplet Selection*: Since the majority of possible triplets trivially fulfill the constraint in (3), improper triplets generating strategies may slow or even interrupt the convergence process during the training. Therefore, it is crucial to select proper triplets. We adopt the online triplet generating strategy [36], i.e., the positive/negative exemplars are selected within a mini-batch during the training process. We sample a certain number of images in every scene-level category in a mini-batch to ensure a certain number of meaningful anchor-positive pairs. In addition, the negative images are randomly sampled in the mini-batch.

Following the study in [36], we use all anchor-positive pairs in a mini-batch and select the semihard negatives, \mathbf{X}_i^n , such that

$$\|f(\mathbf{X}_i^a) - f(\mathbf{X}_i^p)\|_2^2 < \|f(\mathbf{X}_i^a) - f(\mathbf{X}_i^n)\|_2^2. \quad (5)$$

E. Joint Learning

To train our proposed method in an end-to-end fashion, we combine the classification loss and the triplet loss linearly as follows:

$$L = L_c + \lambda L_t \quad (6)$$

where λ is a tradeoff parameter.

IV. EXPERIMENTAL RESULTS AND DISCUSSION

Three RS image data sets are used to evaluate the performance of the proposed method. The first is the multilabel UC Merced Land Use data set¹ (UC-Merced data set) [4], an extension from the UC-Merced data set [43]. The second is the Ankara Hyperspectral Image data set² (Ankara data

set) [44]. The third is the multilabel version of the Aerial Image Data set (AID) data set³ [5], which is extended from the widely used AID data set [16]. Each image in the three data sets has a single label at the scene level and multiple labels at the object level.

First, we perform a group of experiments to evaluate the contributions of different components of our method. Second, comparisons are conducted between the proposed method and state-of-the-art multilabel annotation methods. Then, we show some annotation case results. After that, the model and computational complexity of the proposed method and other deep learning methods are analyzed. Finally, the parameters of the proposed method are analyzed.

A. Data Sets

1) *UC-Merced Data Set*: The UC-Merced data set [43] contains 2100 images that are grouped into 21 broad categories at the scene level. There are 100 images per category with a size of $256 \times 256 \times 3$ and a spatial resolution of 0.3 m. Considering the multilabel property of the images, Chauduri *et al.* [4] relabeled the data set with a redesigned label set. The total number of distinct object-level labels is 17, representing a diverse set of primitive classes in the images: *airplane, bare soil, buildings, cars, chaparral, court, dock, field, grass, mobile home, pavement, sand, sea, ship, tanks, trees, and water*. Each image in the UC-Merced data set is manually annotated with one or more (maximum 7) labels at the object level. On average, each image contains 3.3 object-level labels. Detailed information about this data set can be found in [4] and [43]. Fig. 3 illustrates some example images with single and multiple labels at dual levels. Fig. 6(a) lists the number of images associated with each object-level label.

2) *Ankara Data Set*: The Ankara data set [44] is a small hyperspectral image archive consisting of 216 image tiles with a size of 63×63 pixels. The image tiles are obtained by fragmenting the large hyperspectral images, which were acquired by an EO-1 Hyperion sensor from the area surrounding the city of Ankara in Turkey. The ground resolution is 30 m. Each image tile is associated with multiple object-level labels (land-cover classes) and a single land-use scene-level label. The image tiles are grouped into four categories at the scene level and 29 classes at the object level. On average, each image contains nine object-level labels, and the maximum is 17. Detailed information about this archive can be found in [44]. The Ankara data set contains 119 channels hyperspectral images and the corresponding three channel (RGB) images. We use the RGB images in our experiment. Fig. 4 shows some example images per scene-level category and the associated multiple object-level labels. Fig. 6(b) lists the number of images per object-level label.

3) *Multilabel AID Data Set*: In order to further evaluate our method, we test it and its counterparts on a more challenging data set, the multilabel AID data set [5]. Extended from the widely used AID data set [16], the multilabel AID data set includes 3000 aerial images from 30 categories with manually labeled multiple object labels. With a size of 600×600 pixels,

¹The UC-Merced data set is downloaded from http://bigearth.eu/data_sets

²The Ankara data set is downloaded from http://bigearth.eu/data_sets

³The multilabel AID data set is downloaded from https://github.com/Hua-YIS/AID-Multilabel-Data_set



Fig. 3. Example images from the UC-Merced data set (best viewed in color). The scene-level semantic concepts are shown in black and the object-level semantic concepts are shown in red.

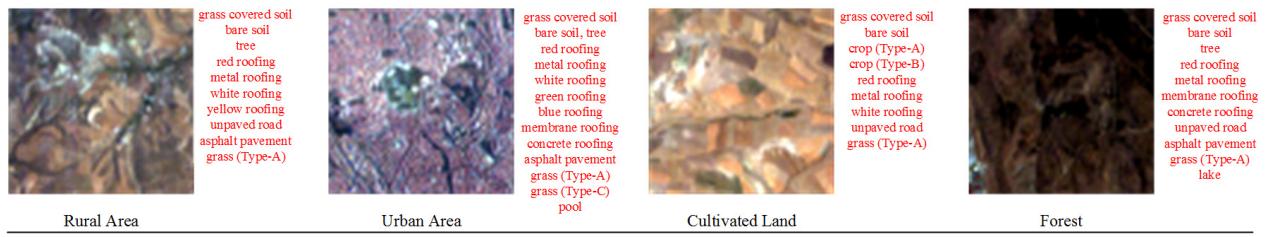


Fig. 4. Example images from the Ankara data set (best viewed in color). The scene-level semantic concepts are shown in black and the object-level semantic concepts are shown in red.

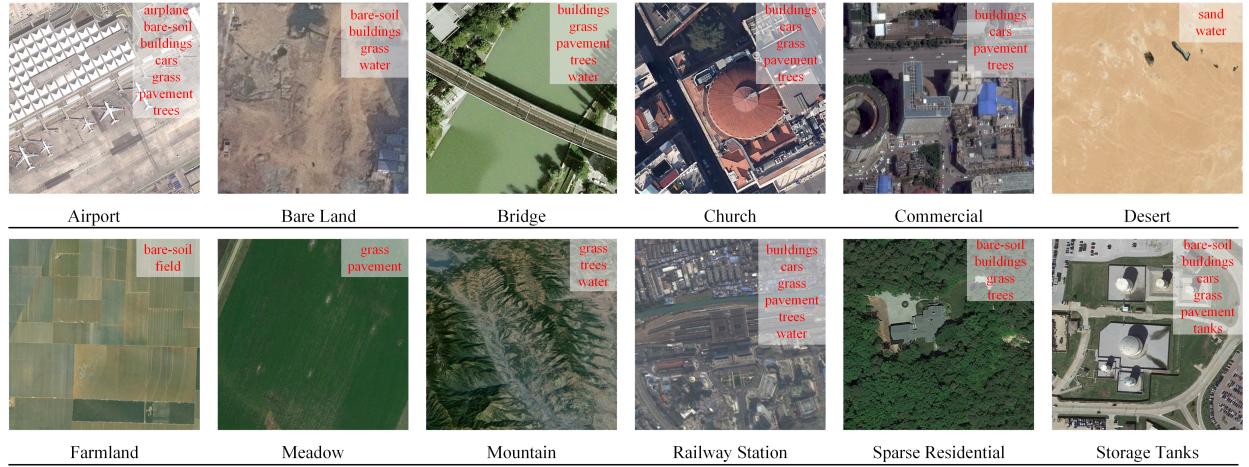


Fig. 5. Example images from the multilabel AID data set (best viewed in color). The scene-level semantic concepts are shown in black and the object-level semantic concepts are shown in red.

each image has an average of 5.5 object-level labels and a maximum of 11 object-level labels. The spatial resolution of images in the AID multilabel data set varies from 0.3 to 8 m. Detailed information about this data set can be found in [5] and [16]. Fig. 5 shows some example images with single and multiple labels at dual levels. Fig. 6(c) lists the number of images associated with each object-level label.

B. Metrics

Existing multilabel metrics, which are mainly divided into two groups, i.e., example-based metrics

and label-based metrics, are used to evaluate the annotation performance from various aspects. In our experiments, nine different metrics are selected to evaluate the performance of the proposed method in comparison with several state-of-the-art methods. First, according to [1], the widely adopted example-based and label-based metrics, including precision, recall, and F1 score, are computed.

As described in [1], the label-based precision (T-P) and recall (T-R) metrics refer to the mean precision and recall over all multiple label classes, respectively. The example-based precision (I-P) and recall (I-R) metrics refer to the aver-

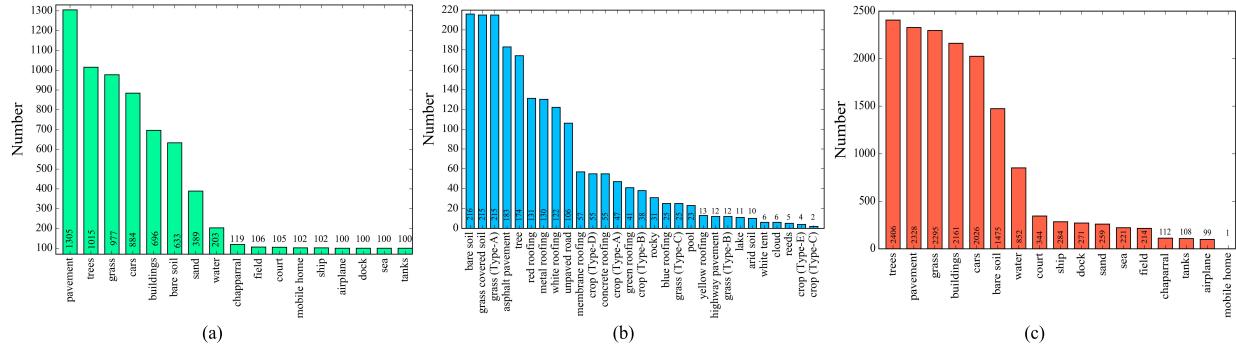


Fig. 6. Number of images per object-level label in the three data sets. (a) UC-Merced data set. (b) Ankara data set. (c) Multilabel AID data set.

age precision and recall over all test images, respectively. The label-based F1 score (T-F1) and example-based F1-score (I-F1) are used as comprehensive performance evaluations by combining the precision and recall with the harmonic mean. The abovementioned six metrics are defined as follows:

$$\begin{aligned} \text{T-P} &= \frac{1}{q} \sum_{j=1}^q \frac{|I_j^p \cap I_j^g|}{|I_j^p|}, \quad \text{I-P} = \frac{1}{n} \sum_{i=1}^n \frac{|T_i^p \cap T_i^g|}{|T_i^p|} \\ \text{T-R} &= \frac{1}{q} \sum_{j=1}^q \frac{|I_j^p \cap I_j^g|}{|I_j^g|}, \quad \text{I-R} = \frac{1}{n} \sum_{i=1}^n \frac{|T_i^p \cap T_i^g|}{|T_i^g|} \\ \text{T-F1} &= \frac{2 \cdot \text{T-P} \cdot \text{T-R}}{\text{T-P} + \text{T-R}}, \quad \text{I-F1} = \frac{2 \cdot \text{I-P} \cdot \text{I-R}}{\text{I-P} + \text{I-R}}. \end{aligned} \quad (7)$$

Here, q stands for the number of multiple labels and n is the number of test images. I_j^p is the set of images that are correctly labeled as class j and I_j^g is the set of images that are associated with a ground-truth label of class j . T_i^p and T_i^g are the predicted and ground-truth label set for image X_i , respectively.

The label-based metrics are biased toward infrequent labels, and the example-based metrics are biased toward frequent labels. Inspired by [23], we define a new metric H-F1, which is the harmonic mean of T-F1 and I-F1

$$\text{H-F1} = \frac{2 \cdot \text{T-F1} \cdot \text{I-F1}}{\text{T-F1} + \text{I-F1}}. \quad (8)$$

Furthermore, the Exact Match and the Hamming Loss metrics are also computed in our experiments. The Exact Match is the ratio of the number of images for which all the labels are correctly predicted, i.e., the predicted label set is identical to the ground-truth label set, to the total number of participating images. The Hamming Loss quantifies the fraction of misclassified image-label pairs, that is, a relevant label is missed, or an irrelevant label is predicted. The Hamming Loss is defined as follows:

$$hloss = \frac{1}{n} \sum_{i=1}^n \frac{1}{q} |T_i^p \Delta T_i^g| \quad (9)$$

where Δ stands for the symmetric difference between two sets. q stands for the number of multiple labels and n is the number of test images.

For each label, we compute the label-specific F1 score, which is based on the label-specific precision and recall.

The F1 score is computed as follows:

$$\text{Precision} = \frac{|I_j^p \cap I_j^g|}{|I_j^p|}, \quad \text{Recall} = \frac{|I_j^p \cap I_j^g|}{|I_j^g|}$$

$$\text{F1} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (10)$$

where I_j^p is the set of images that are correctly labeled as class j and I_j^g is the set of images that are associated with a ground-truth label of class j .

C. Implementation Details

In our experiments, multiple labels are encoded into multihot binary sequences with their length equal to the number of all candidate object labels. In each sequence, 1 indicates the existence of the corresponding object, while 0 denotes the absence. The pretrained weights from ImageNet [45] are used to initialize the shared CNN component of our model and weights of other components are initialized with a Xavier uniform initializer [46]. The training is conducted with the Adam optimizer [47] and the learning rate is initially set to 1e-5, which decays every ten epochs with a factor of 0.9. Other parameters of the optimizer are set as recommended: $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The model is based on Tensorflow and Keras code and is trained on a 1080 Ti GPU with 16-GB memory. We set the batch size to 24 and train the model for 60 epochs. The dimension of the embedding is set to 256. During the training, for each epoch, we first select a certain number of images from each scene category and get their embeddings. Then we find the potential triplets through our triplet selection strategy, rearrange the embeddings in the form of triplets and calculate the loss in batches. In addition, an exhaustive threshold search with a 0.01 increment from 0 to 1 is performed to find class-specific best thresholds. During the test, the output sequences are binarized with the class-specific best thresholds to obtain the final predictions.

D. Evaluation for Model Components

The proposed method can be divided into three closely related methods for component effectiveness evaluation: 1) CNN denotes the original Inception-ResNet-v2 model with the binary cross-entropy loss; 2) CNN + triplet loss refers to the method shown in Fig. 2 with the attention module

TABLE II
COMPONENT EFFECTIVENESS EVALUATION OF THE PROPOSED METHOD ON THE UC-MERCED DATA SET

Level	Models	T-P (%)	T-R (%)	T-F1 (%)	I-P (%)	I-R (%)	I-F1 (%)	H-F1 (%)
single level	CNN	83.10 ± 0.73	81.94 ± 0.72	82.05 ± 0.63	84.27 ± 1.49	86.45 ± 1.08	83.63 ± 1.06	82.83 ± 0.69
dual level	CNN+triplet loss	<u>92.01 ± 0.90</u>	<u>92.46 ± 0.78</u>	<u>92.08 ± 0.79</u>	<u>90.63 ± 0.81</u>	91.68 ± 1.37	<u>90.12 ± 1.01</u>	<u>91.09 ± 0.88</u>
single level	CNN+attention+skip-layer connection	88.91 ± 2.36	87.58 ± 2.42	87.91 ± 2.29	87.93 ± 2.01	88.95 ± 1.50	87.01 ± 1.91	87.45 ± 2.05
dual level	Proposed Method	92.96 ± 0.98	92.60 ± 0.52	92.66 ± 0.47	91.75 ± 0.83	<u>91.65 ± 0.76</u>	90.62 ± 0.62	91.63 ± 0.52

The multi-label annotation results ($AC\% \pm STD\%$) are reported. The best results are shown in bold and the second best results are underlined. For all metrics, the higher the value, the better the evaluation.

TABLE III
COMPONENT EFFECTIVENESS EVALUATION OF THE PROPOSED METHOD ON THE ANKARA DATA SET

Level	Models	T-P (%)	T-R (%)	T-F1 (%)	I-P (%)	I-R (%)	I-F1 (%)	H-F1 (%)
single level	CNN	47.33 ± 3.25	46.23 ± 3.94	43.91 ± 3.44	76.93 ± 4.28	76.72 ± 3.98	75.09 ± 2.89	55.38 ± 3.41
dual level	CNN+triplet loss	<u>53.28 ± 2.66</u>	<u>55.73 ± 1.55</u>	<u>51.80 ± 1.43</u>	<u>82.25 ± 2.49</u>	<u>81.23 ± 2.92</u>	<u>79.65 ± 1.68</u>	<u>62.76 ± 1.31</u>
single level	CNN+attention+skip-layer connection	53.27 ± 3.46	<u>56.78 ± 2.84</u>	<u>52.89 ± 3.17</u>	82.98 ± 1.25	82.76 ± 2.19	80.95 ± 1.01	<u>63.94 ± 2.53</u>
dual level	Proposed Method	54.88 ± 2.13	58.96 ± 4.20	54.66 ± 2.55	81.22 ± 2.46	<u>82.12 ± 2.74</u>	79.60 ± 1.47	64.79 ± 2.01

The multi-label annotation results ($AC\% \pm STD\%$) are reported. The best results are shown in bold and the second best results are underlined. For all metrics, the higher the value, the better the evaluation.

TABLE IV
H-F1 ($AC\% \pm STD\%$) ON THE UC-MERCED DATA SET OF OUR METHOD WITH DIFFERENT ATTENTION MODULES

Attention Module	H-F1 (%)
Non-local block	88.66 ± 0.51
Squeeze-and-Excitation block	<u>91.13 ± 0.39</u>
Trunk-and-mask block	91.24 ± 0.39

The best result is shown in bold and the second best result is underlined.

and the skip-layer connection removed; and 3) the CNN + attention + skip-layer connection denotes the model without the embedding branch. To evaluate the effectiveness of the components, we run each component method and the entire method on the UC-Merced data set and the Ankara data set.

In each data set, we randomly sample 80% of the data as the training samples, 10% of the data as the validation samples, and the remaining 10% of the data as the test samples. We run each experiment over ten random splits and report the mean annotation accuracy (ac) and standard deviation (STD). The results are shown in Tables II and III. From the tables, we make the following observations.

- 1) We report the annotation results of each method utilizing the single- and dual-level semantic concepts. The results show that the utilization of dual-level semantic concepts significantly improves the multilabel annotation performance. For example, the H-F1 score of the CNN + triplet loss method increases 8.26% on the UC-Merced data set and 7.38% on the Ankara data set compared to the CNN method.

2) The effectiveness of the attention module and skip-layer connection is evaluated by comparing the CNN method with and the CNN + attention + skip-layer connection method. The H-F1 scores show a significant increment with the inclusion of the attention module (4.62% on the UC-Merced data set and 8.56% on the Ankara data set).

3) The entire method achieves the best annotation performance in terms of the H-F1 score on the two data sets. Compared with the CNN + triplet loss method on the UC-Merced data set (see Table II) and the CNN + attention + skip-layer connection method on the Ankara data set (see Table III), only small gains (<1% for H-F1) are obtained by the entire framework. However, the performance is still impressive since the entire method is more robust than the CNN + triplet loss method and CNN + attention + skip-layer connection method.

To validate the effectiveness of the trunk-and-mask structure attention module, we conduct experiments on the UC-Merced data set by replacing the trunk-and-mask block with two powerful attention blocks, i.e., the nonlocal block [32] and the squeeze-and-excitation block [33]. To ensure a fair comparison, only one block is added to the last block of the Inception-ResNet-v2-C unit for the three attention blocks.

As indicated in Table IV, the trunk-and-mask block shows the best performance. The trunk-and-mask block has two advantages. First, the structure can emphasize good properties of original features and at the same time gives them the ability to bypass the mask subbranch. Thereby, the feature selection ability of the mask branch is weakened. Second, the

TABLE V
MULTILABEL ANNOTATION RESULTS ($AC\% \pm STD\%$) OF THE PROPOSED METHOD ON THE UC-MERCED DATA SET IN COMPARISON WITH STATE-OF-THE-ART METHODS

Models	T-P (%)	T-R (%)	T-F1 (%)	I-P (%)	I-R (%)	I-F1 (%)	H-F1 (%)	Exact Match (%)	Hamming Loss
ML-KNN	88.63 ± 1.16	84.36 ± 0.90	85.96 ± 0.59	87.28 ± 1.22	86.73 ± 1.27	85.53 ± 0.72	85.74 ± 0.59	48.67 ± 1.20	0.06 ± 0.00
FastTag	84.54 ± 1.85	84.53 ± 1.46	84.00 ± 1.47	78.12 ± 2.87	83.08 ± 2.73	78.20 ± 2.63	80.99 ± 2.03	33.00 ± 3.93	0.08 ± 0.01
CNN+Binary cross-entropy	89.01 ± 1.28	89.05 ± 1.61	88.76 ± 0.87	86.14 ± 1.43	86.04 ± 1.65	84.47 ± 0.91	86.56 ± 0.80	39.00 ± 2.37	0.07 ± 0.00
CNN-RNN	65.70 ± 4.23	62.47 ± 4.15	61.22 ± 3.80	74.79 ± 2.50	79.88 ± 2.60	75.12 ± 2.20	67.44 ± 3.17	30.43 ± 4.10	0.11 ± 0.01
Stivaktakis [49]	9.50 ± 0.78	18.16 ± 1.79	12.41 ± 1.02	52.19 ± 2.56	46.42 ± 3.37	46.62 ± 1.84	19.59 ± 1.38	0.86 ± 1.04	0.19 ± 0.01
Seymour [24]	23.73 ± 1.94	29.54 ± 2.81	24.56 ± 1.22	55.47 ± 3.55	60.16 ± 4.57	56.60 ± 3.50	34.23 ± 1.61	0.29 ± 0.61	0.25 ± 0.02
Proposed Method	92.96 ± 0.98	92.60 ± 0.52	92.66 ± 0.47	91.75 ± 0.83	91.65 ± 0.76	90.62 ± 0.62	91.63 ± 0.52	54.52 ± 2.87	0.04 ± 0.00

The best results are shown in bold and the second best results are underlined. For all metrics, except the Hamming Loss, the higher the value, the better the evaluation. For the Hamming Loss, the lower the value, the better the evaluation.

TABLE VI
MULTILABEL ANNOTATION RESULTS ($AC\% \pm STD\%$) OF THE PROPOSED METHOD ON THE ANKARA DATA SET IN COMPARISON WITH STATE-OF-THE-ART METHODS

Models	T-P (%)	T-R (%)	T-F1 (%)	I-P (%)	I-R (%)	I-F1 (%)	H-F1 (%)	Exact Match (%)	Hamming Loss
ML-KNN	37.14 ± 1.03	36.24 ± 1.25	35.05 ± 1.20	83.84 ± 2.07	73.94 ± 0.66	76.83 ± 1.03	48.13 ± 1.14	0.00 ± 0.00	0.13 ± 0.01
FastTag	45.39 ± 3.76	43.91 ± 4.05	41.39 ± 3.36	82.00 ± 2.45	67.64 ± 3.50	71.10 ± 2.42	52.27 ± 3.14	1.43 ± 0.02	0.16 ± 0.01
CNN+Binary cross-entropy	42.90 ± 4.53	42.93 ± 3.93	41.64 ± 3.94	78.59 ± 4.31	77.73 ± 4.10	76.03 ± 3.81	53.72 ± 3.76	1.36 ± 0.03	0.15 ± 0.02
CNN-RNN	36.25 ± 1.61	50.79 ± 2.73	40.17 ± 1.72	72.59 ± 2.78	85.83 ± 1.61	76.74 ± 1.17	52.71 ± 1.51	0.95 ± 0.02	0.16 ± 0.01
Stivaktakis [49]	24.07 ± 1.37	30.81 ± 1.33	26.56 ± 1.27	77.84 ± 4.17	77.56 ± 3.05	75.42 ± 2.39	39.28 ± 1.63	0.04 ± 0.04	0.15 ± 0.01
Seymour [24]	27.36 ± 1.45	49.58 ± 1.60	31.18 ± 1.55	52.31 ± 3.85	82.55 ± 2.96	62.20 ± 3.51	41.53 ± 2.07	0.09 ± 0.04	0.29 ± 0.03
Proposed Method	54.88 ± 2.13	58.96 ± 4.20	54.66 ± 2.55	81.22 ± 2.46	82.12 ± 2.74	79.60 ± 1.47	64.79 ± 2.01	0.00 ± 0.00	0.12 ± 0.01

The best results are shown in bold and the second best results are underlined. For all metrics, except the Hamming Loss, the higher the value, the better the evaluation. For the Hamming Loss, the lower the value, the better the evaluation.

TABLE VII
MULTILABEL ANNOTATION RESULTS ($AC\% \pm STD\%$) OF THE PROPOSED METHOD ON THE MULTILABEL AID DATA SET IN COMPARISON WITH STATE-OF-THE-ART METHODS

Models	T-P (%)	T-R (%)	T-F1 (%)	I-P (%)	I-R (%)	I-F1 (%)	H-F1 (%)	Exact Match (%)	Hamming Loss
ML-KNN	73.28 ± 1.16	59.31 ± 0.95	64.10 ± 0.85	87.11 ± 0.84	83.44 ± 0.43	83.40 ± 0.53	72.49 ± 0.67	30.87 ± 1.50	0.08 ± 0.00
FastTag	65.71 ± 1.36	62.84 ± 0.57	63.51 ± 0.79	81.41 ± 1.24	80.66 ± 1.00	78.81 ± 0.94	70.33 ± 0.73	18.67 ± 0.94	0.11 ± 0.00
CNN+Binary cross-entropy	69.39 ± 1.41	67.15 ± 2.91	66.82 ± 2.57	82.11 ± 2.06	84.99 ± 1.92	80.94 ± 2.11	73.20 ± 2.40	0.00 ± 0.00	0.11 ± 0.01
CNN-RNN	56.90 ± 5.63	55.72 ± 3.64	54.11 ± 3.81	84.06 ± 2.73	85.01 ± 1.44	82.34 ± 1.60	65.26 ± 3.17	24.00 ± 1.63	0.10 ± 0.01
Stivaktakis [49]	24.80 ± 0.62	35.29 ± 0.00	28.96 ± 0.44	70.26 ± 1.77	81.49 ± 2.04	71.74 ± 1.54	41.26 ± 0.70	0.00 ± 0.00	0.16 ± 0.01
Seymour [24]	39.14 ± 2.47	44.73 ± 7.17	35.83 ± 3.05	65.87 ± 5.02	77.27 ± 10.24	69.82 ± 5.52	47.35 ± 3.88	0.20 ± 0.27	0.25 ± 0.04
Proposed Method	80.89 ± 1.84	74.08 ± 3.11	76.50 ± 2.39	89.72 ± 0.44	88.41 ± 0.65	87.49 ± 0.18	81.61 ± 1.37	22.47 ± 18.46	0.07 ± 0.00

The best results are shown in bold and the second best results are underlined. For all metrics, except the Hamming Loss, the higher the value, the better the evaluation. For the Hamming Loss, the lower the value, the better the evaluation.

upsampling operation in the trunk-and-mask block is the premise of multiscale information fusion. Hence, the network can benefit from the global and local information encoding [35]. The performance of the nonlocal block is the poorest among the three attention modules. One possible reason is that the last block of the Inception-ResNet-v2-C unit is insufficient to provide precise global spatial information with a small spatial size (8×8 pixels). The performance of the squeeze-and-excitation block is slightly poorer than the trunk-and-mask block most likely due to the omission of the spatial attention.

E. Comparisons With State-of-the-Art Methods

To validate the effectiveness of our proposed method for multilabel annotation, we compare it with the following image annotation methods.

- 1) *ML-KNN* [14]: The ML-KNN is a lazy learning multilabel annotation method, which is based on statistical information gained from its KNN and the maximum *a posteriori* principle.
- 2) *FastTag* [48]: In the FastTag method, two coregularized linear mappings are utilized to learn two multilabel

classifiers from the image and text modalities simultaneously.

- 3) *CNN + Binary Cross-Entropy* [21]: This is a deep CNN model that uses a sigmoidal activation function in the final layer and binary cross-entropy as the loss function.
- 4) *CNN-RNN* [22]: This method learns a joint image-label embedding to characterize the semantic label dependence and the image-label relevance.
- 5) *Stivaktakis* [49]: Stivaktakis employed a data augmentation technique to improve the performance of a CNN architecture in the multilabel annotation.
- 6) *Seymour* [24]: Seymour proposed a triplet-based dual-CNN framework to embed images and their labels into a joint space.

In each data set, we randomly divide the data into training, validation and test data sets according to the ratios (80%, 10%, and 10%, respectively). We repeat the experiments ten times using the hold-out sample method [50]. The annotation results ($AC\% \pm STD\%$), including the example-based and label-based precisions, recalls and F1 scores, and the Hamming Loss and Exact Match scores, are listed in Tables V–VII.

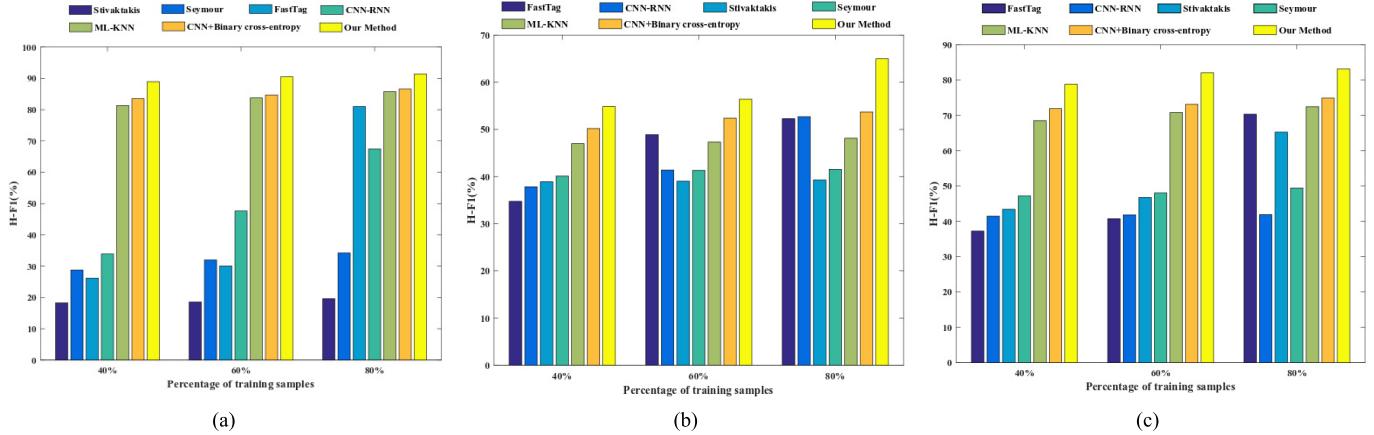


Fig. 7. H-F1 scores on the three data sets of different methods with increasing training ratios. (a) UC-Merced data set. (b) Ankara data set. (c) Multilabel AID data set.

TABLE VIII

PER-CLASS F1 SCORES OF OUR METHOD ON THE UC-MERCED DATA SET IN COMPARISON WITH STATE-OF-THE-ART METHODS

class	Frequency	ML-KNN	FastTag	CNN+Binary cross-entropy	CNN-RNN	Stivaktakis [49]	Seymour [24]	Proposed Method
pavement	0.619	92.90 ± 0.48	89.59 ± 0.88	90.68 ± 2.15	89.18 ± 1.64	76.50 ± 2.33	86.41 ± 2.62	93.29 ± 0.81
trees	0.48	81.91 ± 0.93	77.57 ± 3.55	78.06 ± 2.90	77.70 ± 2.90	64.98 ± 3.75	70.37 ± 12.88	87.86 ± 1.87
grass	0.464	85.80 ± 1.51	78.15 ± 2.36	79.99 ± 3.02	75.87 ± 4.20	63.90 ± 3.08	67.72 ± 10.81	89.13 ± 1.26
cars	0.422	83.95 ± 1.37	78.02 ± 2.32	81.74 ± 2.15	80.81 ± 2.33	5.57 ± 16.70	43.71 ± 17.89	85.52 ± 1.47
bare soil	0.342	74.55 ± 3.93	59.72 ± 2.85	67.30 ± 5.05	51.07 ± 6.17	0.00 ± 0.00	53.61 ± 11.54	78.23 ± 3.64
buildings	0.329	83.20 ± 2.28	80.38 ± 2.87	81.42 ± 3.45	71.95 ± 3.22	0.00 ± 0.00	64.97 ± 10.62	87.98 ± 2.12
sand	0.14	77.59 ± 2.80	75.63 ± 4.32	74.75 ± 4.38	55.23 ± 11.81	0.00 ± 0.00	10.99 ± 9.41	88.86 ± 2.52
water	0.097	87.37 ± 2.50	80.74 ± 3.87	99.28 ± 1.10	66.53 ± 10.29	0.00 ± 0.00	1.93 ± 4.38	97.06 ± 2.36
chaparral	0.055	92.79 ± 2.65	88.77 ± 2.93	86.93 ± 4.86	90.17 ± 6.02	0.00 ± 0.00	1.43 ± 4.29	93.59 ± 1.57
court	0.05	63.25 ± 6.91	67.69 ± 10.08	93.71 ± 3.78	0.00 ± 0.00	0.00 ± 0.00	16.33 ± 20.50	87.05 ± 3.99
field	0.049	88.22 ± 3.30	90.52 ± 1.57	92.32 ± 5.87	82.99 ± 6.70	0.00 ± 0.00	0.00 ± 0.00	94.37 ± 1.73
mobile home	0.049	81.06 ± 3.04	91.89 ± 4.54	91.68 ± 6.55	24.85 ± 21.27	0.00 ± 0.00	0.00 ± 0.00	100.00 ± 0.00
ship	0.049	95.24 ± 0.00	95.56 ± 3.86	98.26 ± 2.17	87.17 ± 6.18	0.00 ± 0.00	0.00 ± 0.00	98.66 ± 2.89
airplane	0.048	95.71 ± 1.43	98.57 ± 2.18	97.96 ± 2.67	8.16 ± 10.44	0.00 ± 0.00	0.00 ± 0.00	100.00 ± 0.00
dock	0.048	100.00 ± 0.00	100.00 ± 0.00	98.81 ± 2.41	90.42 ± 9.14	0.00 ± 0.00	0.00 ± 0.00	99.52 ± 1.43
sea	0.048	92.41 ± 3.08	91.91 ± 4.30	99.52 ± 1.43	80.35 ± 10.30	0.00 ± 0.00	0.00 ± 0.00	98.70 ± 3.91
tanks	0.048	85.44 ± 3.25	83.27 ± 4.82	96.51 ± 4.69	8.26 ± 11.18	0.00 ± 0.00	0.00 ± 0.00	95.32 ± 2.88

The best results are shown in bold and the second best results are underlined.

From Tables V–VII, we make the following observations.

- 1) On the UC-Merced data set (see Table V), all methods obtain good results except the Stivaktakis method [49] and the Seymour method [24]. The poor performance of the Stivaktakis method is due to the simple architecture of three convolutional layers. As the multilabels of a query image are obtained by retrieving from the embedding space, the performance of the Seymour method is limited by the number of images. The proposed

method outperforms the state-of-the-art methods in all metrics, which suggests that the proposed method can significantly improve the multilabel annotation performance. In particular, our proposed method achieves an H-F1 score that is $\geq 5.07\%$ higher than those of the other methods.

- 2) Table VII describes the multilabel annotation performance of each method on the Ankara data set. The acs of all the methods are not very high due to the small number of images and the limited number of several labels

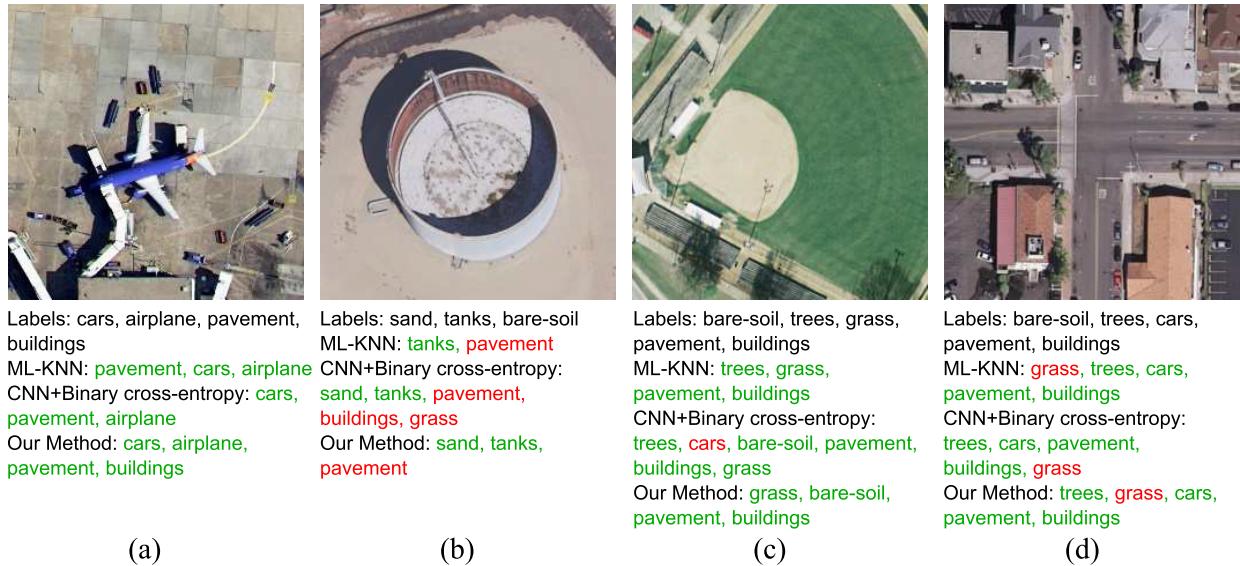


Fig. 8. Annotation example images from the UC-Merced data set with ground-truth object-level labels and predicted object-level labels obtained with ML-KNN, CNN + Binary cross-entropy, and the proposed method. The ground-truth labels are shown in black, the correct predictions are shown in green, and the incorrect predictions are shown in red. (a) Airport. (b) Storage tanks. (c) Baseball diamond. (d) Intersection.

in this data set, e.g., *cloud* and *reeds* [see Fig. 6(b)]. However, the proposed method still outperforms all the other competitors in most of the metrics except the I-P and I-R metrics. This indicates that the proposed method can obtain a better multilabel annotation performance on the hyperspectral data set, even the data set has a limited number of images and labels. Note that the proposed method achieves an H-F1 score that is $\geq 11.07\%$ higher than those of the other methods. This indicates that our proposed method considers both frequent and infrequent labels for multilabel annotation. It is also worth noting that the Exact Match scores of all the methods are zero or close to zero, which is because the Exact Match metric is overly strict and even not applicable in the case of a small data set and a relatively large label space.

- 3) Table VII presents the multilabel annotation performance of each method on the multilabel AID data set. In general, most of the methods exhibit lower performances on the multilabel AID data set than on the UC-Merced data set, which indicates that the multilabel AID data set is more challenging than the UC-Merced data set. Our proposed method achieves the best results on all the metrics except the Exact Match metric. For example, our method obtains a maximum of 81.61% H-F1 score and surpasses the second best method by 8.41%.
- 4) Overall, our proposed approach outperforms the other methods on the RS image multilabel annotation tasks mainly because of the adoption of dual-level semantic concepts. Moreover, the attention mechanism effectively guides the feature learning by focusing on salient objects. Finally, the inclusion of skip-layer connections fuses multiscale information.

We list the corresponding class-specific annotation results ($ac\% \pm STD\%$) on the UC-Merced data set in Table VIII,

which details the F1 scores per object-level label sorted by frequency. Color bars represent the changing trend of the F1 scores with respect to the label frequency. The color gradient from orange to green indicates the change from a small F1 score to a large one.

As shown in Table VIII, our proposed method outperforms the compared methods on most of the labels (12/17 in terms of F1 scores). The advantage of the proposed method can be seen more apparent from the distribution of the color bars. The proposed method achieves more than 80% on the F1 score except for the *bare soil* class, of which the F1 score is 78.23%. The reason may be that the *bare soil* class exhibits a large intraclass heterogeneity. The Stivaktakis method and the Seymour method perform very poorly on classes with low frequency, such as the *field* and *tanks* classes. The frequency of different labels varies greatly. The F1 scores of CNN-RNN, Stivaktakis and Seymour methods are more affected by the frequency than the other methods.

To further validate the effectiveness of our proposed method, we test the performance of our method and the compared methods on the three data sets with different training ratios (40%, 60%, and 80%). The results are shown in Fig. 7. It is noted that the performance of all the methods increases with the increase of the training ratio. However, the performance of some methods improves significantly, such as the FastTag method, while that of some other methods only increases slightly, such as the Stivaktakis method and the Seymour method. Our method outperforms the other methods for all the training ratios.

F. Annotation Case Studies

Fig. 8 displays four example images from the UC-Merced data set, their ground-truth labels at the object level and the annotations obtained using the best three methods. It is observed that the proposed method outperforms ML-KNN

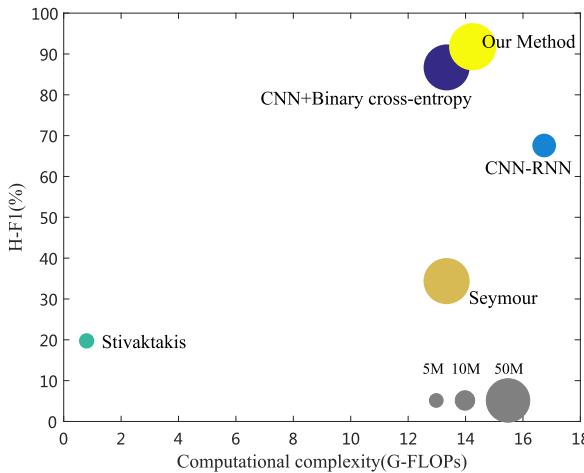


Fig. 9. Ball chart reporting the H-F1 scores of the proposed method and other deep learning methods on the UC-Merced data set versus their computational complexity. Computational complexity is measured by counting FLOPs required for a single forward pass. The ball size corresponds to the number of trainable parameters.

and CNN + Binary cross-entropy on the annotation of the UC-Merced data set. For instance, Fig. 8(a) is annotated correctly by our method while ML-KNN and CNN + binary cross-entropy fail to recall the label *buildings*. Different methods have different preferences when annotating. CNN + Binary cross-entropy tends to recall, whereas ML-KNN tends to precision and our method balances the two metrics. Take Fig. 8(b) for example, CNN + Binary cross-entropy assigns five labels with a precision of 0.40 and a recall of 0.67; ML-KNN assigns two labels with a precision of 0.50 and a recall of 0.33, and our method assigns three labels with a precision of 0.67 and a recall of 0.67.

G. Model and Computational Complexity Analysis

We report the H-F1 scores on the UC-Merced data set versus the computational complexity of our proposed method and other deep learning models in Fig. 9. The size of the balls in Fig. 9 corresponds to the model complexity. We evaluate model complexity by counting the total amount of trainable parameters and measure computational complexity using the floating-point operations (FLOPs) [51] required for a single forward pass. The multiply adds are counted as two FLOPs as was done in [52]. All the considered models expect one input image with $229 \times 229 \times 3$ pixels, except CNN-RNN, which expects $224 \times 224 \times 3$ pixels.

It can be seen from Fig. 9 that our model contains 5.91×10^6 parameters and requires ~ 14.27 G-FLOPs. It is relatively complex but powerful. For example, compared to the second best CNN + binary cross-entropy model, the proposed model achieves a 5.86% increase in H-F1 score with a 6.81% increase in computational complexity and 5.84% increase in model complexity. In addition, on a 1080 Ti GPU with 16 GB memory, the proposed model is capable of super real-time performance (at 60 frames/s) [51] with 10.48-ms inference time per image. Overall, our proposed method can significantly improve the annotation performance with acceptable additional computational and model burden.

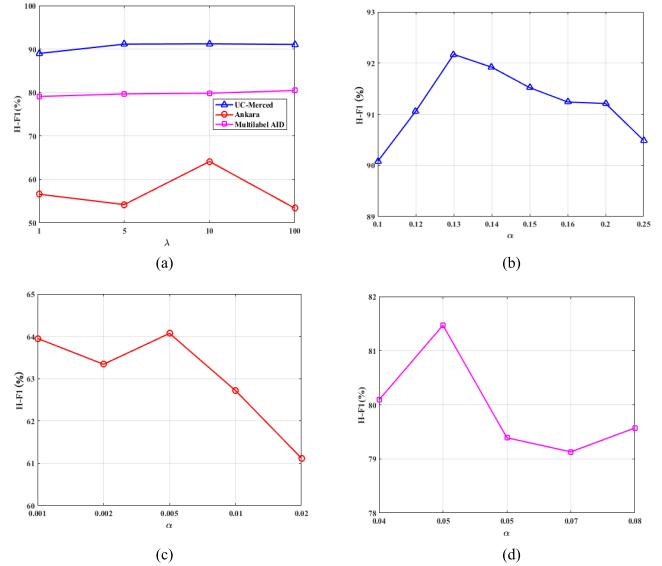


Fig. 10. Influence of parameter selection on the H-F1 scores on the UC-Merced, Ankara, and multilabel AID data sets. (a) H-F1 scores on the UC-Merced, Ankara, and multilabel AID data sets with a varying λ . (b) H-F1 scores on the UC-Merced data set with a varying α . (c) H-F1 scores on the Ankara data set with a varying α . (d) H-F1 scores on the multilabel AID data set with a varying α .

H. Parameter Sensitivity Analysis

In our method, the parameters λ and α need to be tuned manually for each data set. λ is a tradeoff parameter between the classification loss and the triplet loss. α refers to the enforced margin between the positive and negative image pairs when constructing triplets. The results are shown in Fig. 10.

From Fig. 10, it is observed that the optimal parameters are different for the three data sets and the setting of the parameters is important for the annotation performance. With $\lambda = 10$, our method obtains the highest H-F1 score on all the three data sets. Note that our method is not sensitive to the variation of λ on the UC-Merced data set. The H-F1 score peaks at $\alpha = 0.13$, $\alpha = 0.005$, and $\alpha = 0.05$ for the UC-Merced data set, the Ankara data set, and the multilabel AID data set, respectively.

V. CONCLUSION

In this article, we propose an end-to-end deep learning framework for multilabel annotation of RS images that exploits dual-level semantic concepts. The framework includes a shared CNN for visual feature learning, a classification branch for multilabel annotation and an embedding branch for maintaining the similarity relationships among the triplet images grouped by scene-level semantic concepts. An attention mechanism is introduced in the classification branch for salient object detection, while the skip connection is incorporated to combine information from multiple layers. The experimental results of three RS image data sets indicate that the proposed framework leads to significant improvement in the multilabel annotation performance compared to other annotation methods. Our proposed method demonstrates a superior performance with an H-F1 score that is $\geq 5.07\%$ higher on the UC-Merced data set, $\geq 11.07\%$ higher on the Ankara data set

and $\geq 8.41\%$ higher on the multilabel AID data set than the other models.

The Ankara data set provides 119 channels hyperspectral images and the corresponding three channel (RGB) images. In the current study, we only used the three channels (RGB) images. The main reason is for the fairness of comparison. The compared methods in the experiments of the RM utilize either three channel (RGB) images or the features extracted from RGB images. To directly handle images with 119 channels, we may add an adaptation layer to the beginning of our network. Then the model uses 1×1 convolution to reduce the channels and stack Inception-ResNet-v2 on the top of the adaptation layer. In this way, our model can use the available information of hyperspectral images. To make a tradeoff between accuracy versus computational performance depending on the number of channels used, we also can use PCA to reduce the spectral dimension, extract some principal components and send them to our model. In the future work, we continue to improve the study.

The main drawback of the proposed method is that it fails to consider the label dependences at the object level and the label relationships between the scene level and the object level. Thus, in the future work, we plan to adopt the RNN to model the label relationships between the intralevel and interlevel semantic concepts of RS images. In addition, we fail to take advantage of the image embeddings generated in our proposed method. Therefore, as future development of this article, we will use the generated image embeddings for the retrieval tasks and form a multitask learning framework.

REFERENCES

- [1] M.-L. Zhang and Z.-H. Zhou, "A review on multi-label learning algorithms," *IEEE Trans. Knowl. Data Eng.*, vol. 26, no. 8, pp. 1819–1837, Aug. 2014.
- [2] Y. Wang *et al.*, "A three-layered graph-based learning approach for remote sensing image retrieval," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 10, pp. 6020–6034, Oct. 2016.
- [3] Y. Li, Y. Zhang, X. Huang, H. Zhu, and J. Ma, "Large-scale remote sensing image retrieval by deep hashing neural networks," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 2, pp. 950–965, Feb. 2018.
- [4] B. Chaudhuri, B. Demir, S. Chaudhuri, and L. Bruzzone, "Multilabel remote sensing image retrieval using a semisupervised graph-theoretic method," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 2, pp. 1144–1158, Feb. 2018.
- [5] Y. Hua, L. Mou, and X. X. Zhu, "Relation network for multi-label aerial image classification," 2019, *arXiv: 1907.07274*. [Online]. Available: <https://arxiv.org/abs/1907.07274>
- [6] Y. Hua, L. Mou, and X. X. Zhu, "Recurrently exploring class-wise attention in a hybrid convolutional and bidirectional LSTM network for multi-label aerial image classification," *ISPRS J. Photogram. Remote Sens.*, vol. 149, pp. 188–199, Mar. 2019.
- [7] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. CVPR*, Jun. 2015, pp. 3431–3440.
- [8] Z. Zhang, Q. Liu, and Y. Wang, "Road extraction by deep residual U-net," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 5, pp. 749–753, May 2018.
- [9] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Proc. MICCAI*, 2015, pp. 234–241.
- [10] K. Bittner, F. Adam, S. Cui, M. Körner, and P. Reinartz, "Building footprint extraction from VHR remote sensing images combined with normalized DSMs using fused fully convolutional networks," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 11, no. 8, pp. 2615–2629, Aug. 2018.
- [11] K. Karalas, G. Tsagkatakis, M. Zervakis, and P. Tsakalides, "Land classification using remotely sensed data: Going multilabel," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 6, pp. 3548–3563, Jun. 2016.
- [12] W. Boonpook, Y. Tan, Y. Ye, P. Torteeka, K. Torsri, and S. Dong, "A deep learning approach on building detection from unmanned aerial vehicle-based images in riverbank monitoring," *Sensors*, vol. 18, no. 11, p. 3921, 2018.
- [13] G. Tsoumacas, I. Katakos, and I. Vlahavas, "Mining multi-label data," in *Data Mining Knowledge Discovery Handbook*. Boston, MA, USA: Springer, 2009, pp. 667–685.
- [14] M.-L. Zhang and Z.-H. Zhou, "ML-KNN: A lazy learning approach to multi-label learning," *Pattern Recognit.*, vol. 40, no. 7, pp. 2038–2048, 2007.
- [15] Q. Tan, Y. Liu, X. Chen, and G. Yu, "Multi-Label classification based on low rank representation for image annotation," *Remote Sens.*, vol. 9, no. 2, p. 109, Jan. 2017.
- [16] G.-S. Xia *et al.*, "AID: A benchmark data set for performance evaluation of aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 7, pp. 3965–3981, Jul. 2017.
- [17] N. He, L. Fang, S. Li, A. Plaza, and J. Plaza, "Remote sensing scene classification using multilayer stacked covariance pooling," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 12, pp. 6899–6910, Dec. 2018.
- [18] L. Mou, P. Ghamisi, and X. X. Zhu, "Deep recurrent neural networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 7, pp. 3639–3655, Jul. 2017.
- [19] P. Zhu *et al.*, "Projection learning with local and global consistency constraints for scene classification," *ISPRS-J. Photogramm. Remote Sens.*, vol. 144, pp. 202–216, Oct. 2018.
- [20] E. Maggiore, Y. Tarabalka, G. Charpiat, and P. Alliez, "High-resolution aerial image labeling with convolutional neural networks," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 12, pp. 7092–7103, Dec. 2017.
- [21] D. Gardner and D. Nichols. (2017). *Multi-Label Classification of Satellite Images With Deep Learning*. [Online]. Available: <http://cs231n.stanford.edu/reports/2017/pdfs/908.pdf>
- [22] J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, and W. Xu, "CNN-RNN: A unified framework for multi-label image classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 2285–2294.
- [23] Y. Niu, Z. Lu, J.-R. Wen, T. Xiang, and S.-F. Chang, "Multi-modal multi-scale deep learning for large-scale image annotation," *IEEE Trans. Image Process.*, vol. 28, no. 4, pp. 1720–1731, Apr. 2019.
- [24] Z. Seymour and Z. M. Zhang, "Multi-label triplet embeddings for image annotation from user-generated tags," in *Proc. ACM ICMR*, 2018, pp. 249–256.
- [25] H. Blockeel, L. Schietgat, J. Struyf, S. Džeroski, and A. Clare, "Decision trees for hierarchical multilabel classification: A case study in functional Genomics," in *Proc. PKDD*, 2006, pp. 18–29.
- [26] H. Hu, G.-T. Zhou, Z. Deng, Z. Liao, and G. Mori, "Learning structured inference neural networks with label relations," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 2960–2968.
- [27] V. Ranjan, N. Raswiasia, and C. Jawahar, "Multi-label cross-modal retrieval," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 4094–4102.
- [28] F. Hu, W. Yang, J. Chen, and H. Sun, "Tile-level annotation of satellite images using multi-level max-margin discriminative random field," *Remote Sens.*, vol. 5, no. 5, pp. 2275–2291, 2013.
- [29] C. O. Dumitru, G. Schwarz, and M. Datcu, "Land cover semantic annotation derived from high-resolution SAR images," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 9, no. 6, pp. 2215–2232, Jun. 2016.
- [30] X. Zhang, S. Du, and Q. Wang, "Hierarchical semantic cognition for urban functional zones with VHR satellite images and POI data," *ISPRS J. Photogram. Remote Sens.*, vol. 132, pp. 170–184, Oct. 2017.
- [31] A. Vaswani *et al.*, "Attention is all you need," in *Proc. NIPS*, 2017, pp. 5998–6008.
- [32] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018, pp. 7794–7803.
- [33] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018, pp. 7132–7141.
- [34] X. Zhang, T. Wang, J. Qi, H. Lu, and G. Wang, "Progressive attention guided recurrent network for salient object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018, pp. 714–722.
- [35] F. Wang *et al.*, "Residual attention network for image classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 3156–3164.

- [36] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 815–823.
- [37] K. Q. Weinberger, J. Blitzer, and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," in *Proc. NIPS*, 2005, pp. 1473–1480.
- [38] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng, "Person re-identification by multi-channel parts-based CNN with improved triplet loss function," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 1335–1344.
- [39] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Proc. AAAI*, 2017, pp. 4278–4284.
- [40] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 2818–2826.
- [41] C. M. Bishop, *Pattern Recognition and Machine Learning* (Information Science and Statistics). New York, NY, USA: Springer-Verlag, 2006. [Online]. Available: <http://users.isr.ist.utl.pt/~wurmd/Livros/school/Bishop%20-%20Pattern%20Recognition%20And%20Machine%20Learning%20-%20Springer%20%202006.pdf>
- [42] F. Huang, X. Zhang, Z. Li, T. Mei, Y. He, and Z. Zhao, "Learning social image embedding with deep multimodal attention networks," in *Proc. ACM MM*, 2017, pp. 460–468.
- [43] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. ACM SIGSPATIAL*, 2010, pp. 270–279.
- [44] F. Ömrüuzun, B. Demir, L. Bruzzone, and Y. Y. Çetin, "Content based hyperspectral image retrieval using bag of endmembers image descriptors," in *Proc. 8th Workshop Hyperspectral Image Signal Process., Evol. Remote Sens. (WHISPERS)*, Aug. 2016, pp. 1–4.
- [45] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 248–255.
- [46] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proc. 13th Int. Conf. Artif. Intell. Statist.*, Mar. 2010, pp. 249–256. [Online]. Available: <http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>
- [47] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv: 1412.6980*. [Online]. Available: <https://arxiv.org/abs/1412.6980>
- [48] M. Chen, A. Zheng, and K. Weinberger, "Fast image tagging," in *Proc. ICML*, 2013, pp. 1274–1282.
- [49] R. Stivaktakis, G. Tsagkatakis, and P. Tsakalides, "Deep learning for multilabel land cover scene categorization using data augmentation," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 7, pp. 1031–1035, Jul. 2019.
- [50] T. Mitchell, B. Buchanan, G. DeJong, T. Dietterich, P. Rosenbloom, and A. Waibel, "Machine learning," *Annu. Rev. Comput. Sci.*, vol. 4, no. 1, pp. 417–433, 1990.
- [51] S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64270–64277, 2018.
- [52] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 1492–1500.



Panpan Zhu is currently pursuing the Ph.D. degree with the Faculty of Geographical Science, Beijing Normal University, Beijing, China.

Her research interests include remote sensing image processing and image-based classification and retrieval.



Yumin Tan received the Ph.D. degree in cartography and geographic information system (GIS) from the Institute of Geosciences and Resources, Chinese Academy of Science, Beijing, China, in 2004.

She is currently an Associate Professor with the School of Transportation Science and Engineering, Beihang University, Beijing. Her research interests include remote sensing and GIS space technology application, high-resolution remote sensing information extraction, and point cloud data processing.



Liqiang Zhang received the Ph.D. degree in geoinformatics from the Institute of Remote Sensing Applications, Chinese Academy of Sciences, Beijing, China, in 2004.

He is currently a Professor with the Faculty of Geographical Science, Beijing Normal University, Beijing. His research interests include remote sensing image processing, 3-D urban reconstruction, and spatial object recognition.



Yuebin Wang received the Ph.D. degree from the School of Geography, Beijing Normal University, Beijing, China, in 2016.

He was a Post-Doctoral Researcher with the School of Mathematical Sciences, Beijing Normal University. He is currently an Assistant Professor with the School of Land Science and Technology, China University of Geosciences, Beijing. His research interests include remote sensing imagery processing and 3-D urban modeling.



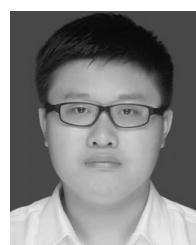
Jie Mei is currently pursuing the Ph.D. degree with the College of Computer Science, Nankai University, Tianjin, China.

His research interests include computer vision, machine learning, and remote sensing image processing.



Hao Liu received the bachelor's degree from Beihang University (BUAA), Beijing, China, in 2012, and the master's degree from Kyushu University, Fukuoka, Japan, in 2015. He is currently pursuing the Ph.D. degree with the Faculty of Geographical Science, Beijing Normal University, Beijing.

His research interests include machine learning and remote sensing image classification.



Mengfan Wu is currently pursuing the Ph.D. degree with the Faculty of Geographical Science, Beijing Normal University, Beijing, China.

His research interests include deep learning and remote sensing image classification.