

An Improved Pretraining Strategy-Based Scene Classification With Deep Learning

Zongli Chen, Yiyue Wang, Wei Han[✉], Ruyi Feng, and Jia Chen[✉]

Abstract—High-resolution remote sensing (HRRS) image scene classification takes an important role in many applications and has attracted much attention. Recently, notable efforts have been made to present massive methods for HRRS scene classification, wherein deep-learning-based methods demonstrate remarkable performance compared with state-of-the-art methods. However, HRRS images contain complex contextual relationships and large differences of object scale, which are significantly different from natural images. The existing deep-learning-based scene classification methods are originally designed for natural image processing and have not been optimized to adapt to the characteristics of HRRS images, which significantly affects the efficiency of the feature extraction and recognition accuracy. In addition, when designing a model for remote sensing tasks, the pretraining of the model is time-consuming. The enormous amount of pretraining time and computation resources necessarily increase the difficulty of producing an excellent model. In this letter, focusing on the problems above, we proposed a new convolutional neural network (CNN)-based scene classification method. The CNN-based scene classification method is constructed by *spatial-scale-aware* blocks and is efficient in extracting the abundant spatial features, but can also adaptively adjust feature responses to maximize the function of informative features in the classification results. In addition, an HRRS imagery-based learning strategy is utilized to obtain an initial model for fine-tuning the model parameters, which drastically reduces the pretraining time. The proposed method has been demonstrated using two HRRS data sets, and experimental results have proven the superiority of the proposed method.

Index Terms—Deep learning, high-resolution remote sensing (HRRS) scene classification, spatial coding, weight-adaptive.

I. INTRODUCTION

THE technologies presently available (e.g., multispectral/hyperspectral, synthetic aperture radar) for earth observation generate many types of aerial and satellite imagers with high resolutions (e.g., spatial resolution, spectral resolution, and temporal resolution). High-resolution remote

sensing (HRRS) image scene classification, which categorizes scene images into an independent set of semantic-level land use and land cover (LULC) classes according to the image contents [1], [2], can extract valuable information from a massive amount of aerial and satellite images. Therefore, it has attracted increasing attention and has been applied in a wide range of applications, such as LULC determination urban planning, environmental protection, and crop monitoring. During the past few decades, many remarkable efforts have been made to develop various methods for HRRS scene classification.

Traditionally, handcrafted-feature-based methods were presented to solve the problem of remote sensing image classification [3]. However, designing handcrafted features consume a remarkable amount of engineering skills and domain expertise. Additionally, the representational capabilities of the handcrafted features are easily influenced by human ingenuity. On the other hand, the deep-feature-based methods, which can make full use of massive annotation data to automatically learn effective features, achieve more impressive results than the handcrafted-feature-based methods in many fields [1], [4]. Therefore, the deep-feature-based methods have been widely applied in natural image classification, object recognition, natural language [5], [6], and have also been used to analyze HRRS images for scene classification [3].

For HRRS scene classification, the widely used convolutional neural network (CNN) models, including visual geometry group (VGG) [5], AlexNet [6], and ResNet [4], are originally designed for natural image processing. However, unlike natural images, HRRS images contain complex contextual relationships and large differences in scales of objects in the images, which significantly affects the efficiency of the feature extraction and recognition accuracy. Many efforts have been made to improve the CNN models for processing remote sensing images [7], [8]. Yu and Liu [9] improved the feature representation by fusing the features from multiple CNN layers. Zhao *et al.* [10] combined spectral and structural information to enhance the representation ability of deep learning features. Some works enhance the performance of the CNN models by improving the spatial coding capability of the models, such as the Inception module [11], which is a multi-branch network [12] for better extraction of spatial features. Some other works focus on the adjustment of feature channel weights, which adaptively recalibrate channelwise feature responses to suppress useless feature channels [13]. The enhancement of the spatial representation and the adaptive recalibration of channelwise feature responses can make the CNN-based models more efficient [13] and more suitable for processing complex HRRS images.

Manuscript received November 25, 2018; revised March 20, 2019 and May 31, 2019; accepted August 6, 2019. This work was supported in part by the National Natural Science Foundation of China under Grant 41571413, Grant 41701429, and Grant U1711266 and in part by the GF Innovative Research Program. (Corresponding authors: Wei Han; Ruyi Feng.)

Z. Chen is with the Department of Land and Resources of Guizhou Province, Guiyang 550004, China.

Y. Wang is with the School of Computer Science, Northeast Forestry University, Harbin 150040, China.

W. Han, R. Feng, and J. Chen are with the School of Computer Science, China University of Geosciences, Wuhan 430074, China, and also with the Hubei Key Laboratory of Intelligent Geo-Information Processing, China University of Geosciences, Wuhan 430074, China (e-mail: weihan@cug.edu.cn; fengry@cug.edu.cn).

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

Digital Object Identifier 10.1109/LGRS.2019.2934341

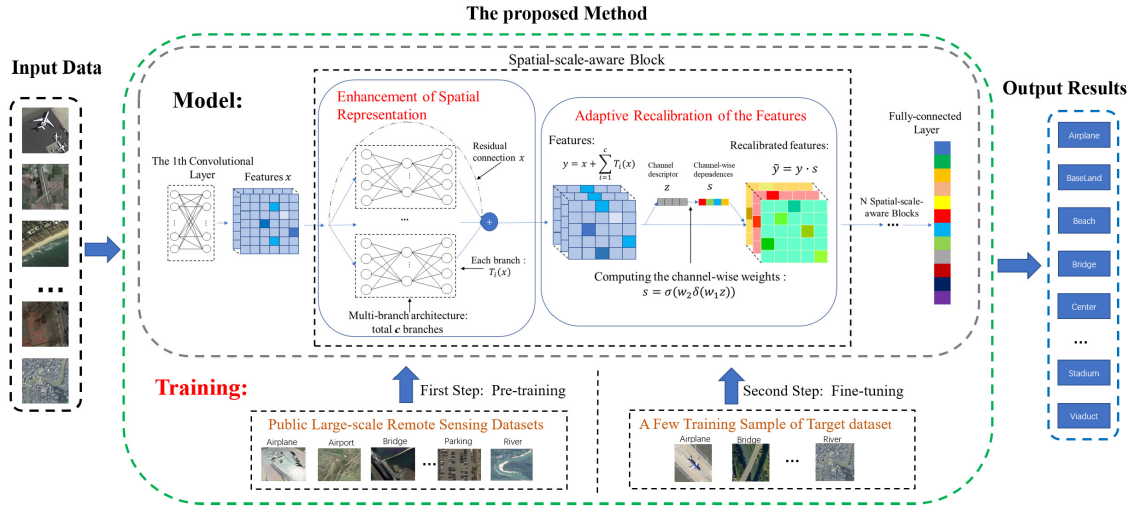


Fig. 1. Process flow of the proposed classification method.

The use of CNN models, pretraining with ImageNet [14] for feature extraction, or fine-tuning of the model parameters has been applied in HRRS scene classification [3]. The pre-training CNN-based models can significantly achieve better classification accuracy and speed up the convergence compared to the CNN models trained from the scratch. However, most of the notable CNN models are developed for natural images. When designing a new model for remote sensing classification, it consumes much time and many computation resources to produce a pretraining model due to the tremendous size of the ImageNet database [14]. Meanwhile, a lot of large-scale HRRS data sets are available in recent years, such as aerial image data set (AID) [2], NWPU-RESISC45 [1], and PatternNet [15]. The massive HRRS data sets make it possible to acquire the models pretraining on HRRS data sets. Because of the similarity of the public HRRS data sets and the target data set, the models pretraining on the HRRS data sets might enhance the classification performance and accelerate the convergence of the CNN models.

Focusing on the challenges of complex contextual relationships, large differences in scale of objects in HRRS images, and the pretraining difficulty of deep learning models, we present a new CNN-based method for HRRS scene classification by utilizing a new convolutional neural block, a *spatial-scale-aware*, and an HRRS imagery-based learning strategy. The spatial-scale-aware block is able to enhance the spatial representation ability and to optimize feature response weights, which is implemented by combining the squeeze and excitation (SE) module [13] and the aggregated transformation (AT) block [12]. In addition, the HRRS imagery-based learning strategy is adapted to finish the model initialization. Compared to the traditional model of pretraining on ImageNet, the proposed strategy is able to speed up and improve the convergence degree of the classification, while significantly reducing pretraining iteration. The experimental results show that the proposed classification model is able to achieve remarkable performance for large-scale HRRS scene classification.

The contributions of this work are summarized as follows.

- 1) A new CNN-based scene classification model is used to achieve the extraction of high-level and multi-scale

features of the objects for HRRS scene classification, wherein the network unit, *spatial-scale-aware block*, can improve the spatial coding, and adaptively adjust the channelwise feature weights.

- 2) An HRRS imagery-based learning strategy is proposed to obtain an initial model for fine-tuning of the model parameters. It can significantly reduce the pretraining iterations, and improve the convergence degree of the proposed network in the target data sets compared with ImageNet pretraining.

The remainder of this letter is organized as follows. In Section II, we introduce the details of the proposed scene classification method. The experiments and analysis are detailed in Section III. Conclusions are drawn in Section IV. The related references are displayed to close.

II. PROPOSED METHOD

To adapt to the characteristics of the complex context and the large differences of object scale in HRRS images, a new HRRS scene classification method is presented. As shown in the main flowchart in Fig. 1, the proposed method contains two main components: the classification network and the learning strategy. The input of the model is a group of HRRS images. The input images are processed by a series of spatial-scale-aware blocks, and high-level feature vectors are obtained. Next, the classification labels of the input images are computed by a following fully connected layer according to the high-level feature vectors. The spatial-scale-aware block plays a critical role in enhancing the discriminability of features, wherein each branch of a multi-branch unit is utilized to extract abundant spatial information in the feature subspace. Then, the channel-correlation weights are calculated by a fully connected layer, which further enhances useful feature channels and suppresses useless feature channels. The learning strategy contains two stages: the pretraining and the fine-tuning. In the pretraining stage, several released HRRS data sets are used to train the classification network and initialize the parameters. Then, in the fine-tuning stage, a few samples of the target data set are used to further adjust the model parameters and to achieve better performance for the proposed model. In Sections II-A and II-B, the more detailed

discussions are provided regarding the classification model and the learning strategy.

A. Classification Network

As mentioned above, the proposed classification consists of a convolutional layer (including nonlinearity and pooling layers), many spatial-scale-aware blocks, and a following fully connected layer for the output of the classification results. The introduction of the spatial-scale-aware blocks, as shown in Fig. 1, are one of the main contributions of the proposed classification method. These play an important role to make the high-level features more discriminative, which enable the extraction of sufficient spatial information by a multi-branch architecture, the calculation of the feature channel correlation, and the adaptive adjustment of the feature channel weights. By adaptive recalibration of the feature responses, informative channels have larger weights and a greater impact on the classification results. In Sections II-A1 and II-A2, the enhancement effects of the spatial representation and the recalibration of feature responses for the spatial-scale-aware block are described in detail.

1) *Enhancement of Spatial Representation*: The first part of the spatial-scale-aware unit is the AT block. It is a homogeneous, multi-branch architecture with a few hyperparameters, and a new dimension, “cardinality,” which is defined as the number of the branch paths in the AT block. Cardinality is different from the existing dimensions of depth and width. The function of the AT block is to recast as a combination of splitting, transforming, and aggregating, which is depicted in Fig. 2. In Fig. 2, the operations of the AT block with a residual connection is formulated as

$$y = x + \sum_{i=0}^C T_i(x) \quad (1)$$

where x symbolizes the feature vectors extracted by the last convolutional layer, $\{T_1, \dots, T_C\}$ are the convolutional operations in each branch of the block, y represents the output features, and C is the value of the cardinality.

The multi-branch structure can map input vectors into different subspaces and to further learning spatial information in the subspace. After that, the features learned in different subspaces are combined as the output features. This architecture is equivalent to the Inception module [16], which is able to learn more spatial information. According to [17], the computational difficulty of deep neural networks is their nonconvex nature. The multi-branch architecture is less nonconvex in terms of the duality gap, which is proven to measure the degree of intrinsic nonconvexity. Therefore, the multi-branch network is more efficacious for reaching a more optimal solution. Increasing the cardinality is also proven to be more effective than increasing the network depth [12].

2) *Recalibration of Feature Responses*: Since the output of a convolutional layer is a summation operation through the feature channels, all of the feature channels are implicit with the same weights. The spatial information learned by each feature channel is different; while some channels might be sensitive to the class of forest, others are greatly responsive to the class of beach. Therefore, if the feature responses can

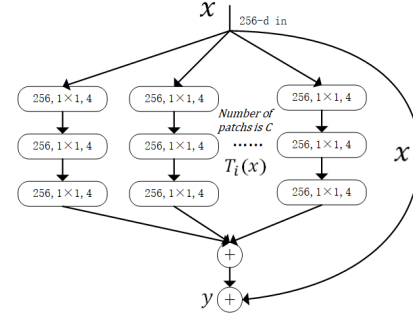


Fig. 2. Illustration of the AT block with residual connection. A block with C branches. A layer is shown as (#in channels, filter size, and #output channels) [12].

be adaptively recalibrated by increasing the weights of the informative channels and suppressing the useless channels, the output features can be more effective in improving the classification performance. The SE module is an available solution to meet the requirement.

In the SE module, a channel descriptor is conducted first by using global average pooling to generate channelwise statistics, compressing global spatial information. Formally, the statistic $z \in R^A$ is generated by shrinking each channel of x through the spatial dimensions $W \times H$. For $a = \{1, \dots, C\}$, x_a is the a th feature response of x and z_a is the a th element of z . z_a is calculated by

$$z_a = \frac{1}{W \times H} \sum_{i=1}^W \sum_{j=1}^H x_a(i, j). \quad (2)$$

Then, to capture the channelwise dependencies s , a gating mechanism with a sigmoid activation σ is utilized

$$s = \sigma(g(z, w)) = \sigma(w_2 \delta(w_1 z)) \quad (3)$$

where δ is the rectified linear unit function. To limit the model's complexity, the SE module is constructed by forming a bottleneck with two fully connected layers for dimensionality reduction around the nonlinearity. The two fully connected layers have the parameters $w_1 \in R^{(C/r) \times C}$ and $w_2 \in R^{C \times (C/r)}$. We set the reduction ratio $r = 16$ to consistent with the original work [13], because this setting achieved a good tradeoff between accuracy and complexity. The final output \tilde{y} is obtained by rescaling the transformation y with the weights s

$$\tilde{y} = y \cdot s. \quad (4)$$

The SE module for ResNet-18 introduces ~ 0.7 million additional parameters beyond the ~ 11.7 million parameters required by ResNet-18, corresponding to a $\sim 6\%$ increase. Therefore, it can be exploited by subsequent transformations and utilized to suppress useless ones. It achieves the function by explicitly modeling channel interdependencies and recalibrating filter responses in two steps, SE, before they are fed into the next convolutional layer.

B. HRRS Imagery-Based Learning Strategy

The learning strategy in the proposed method is discussed here in detail. In the past few years, a well-known paradigm

TABLE I
PROPERTIES OF NWPU-RESISC45, PATTERNNET, AND RSI-CB

Dataset Name	Image/Class	Classes	Size	Images
NWPU-RESISC45[1]	700	45	256×256	31,500
PatternNet [15]	800	38	256×256	30,400
RSI-CB[18]	~690	35	256×256	24,747
UCM[19]	100	21	256×256	2,100
AID[2]	200~400	30	600×600	10,000

of deep learning has been to pretrain deep learning models using ImageNet and then to fine-tune the models on target data set. The paradigm has achieved many state-of-the-art records, such as image classification, object detection, and image segmentation. In the remote sensing community, the pre-training is mainly based on ImageNet [14] due to the lack of annotation HRRS samples. A recent work [20] revealed that ImageNet pretraining speeds up the convergence early in learning process, but not necessarily provide regularization or improve final target task accuracy. Empirically, because of the similarity of the target data and the public HRRS data sets, the use of HRRS data sets for the pretraining can speed up the convergence of the classification model faster than ImageNet pretraining. Thanks to the recently released large-scale HRRS data sets, such as NWPU-RESISC45 [1], PatternNet [15], and RSI-CB [18], the annotated scene images have increased dramatically. The massive number of annotated scene samples make the HRRS imagery-based pretraining strategy possible.

Through the aforementioned analysis, we try to use HRRS for pretraining the CNN models. The total learning process contains two parts: the pretraining stage and fine-tuning stage. In the pretraining stage, the NWPU-RESISC45 [1], PatternNet [15], and RSI-CB [18] data sets, were chosen to build a large-scale HRRS scene data set. The properties of these data sets are shown in Table I. The images with the same class labels in different data sets were merged. The merged HRRS data set contained 86 classes with a total of 85 647 images. The batch size was set to 256, and the number of iterations was 100 000 (about 298 epochs) for pretraining. The learning rate was initially set to 0.1, and the attenuation was 0.1 for every 25 000 iterations. In the fine-tuning stage, the classification model with parameter initialization was trained by a few samples of the target data set. By doing this, the model parameters were further adjusted to fit the target data set. The parameter settings in the stage is the same as a known stable setup [21]. The last fully connected layer of the model was replaced with a new one for adjusting the number of output vectors to the class number of the target data set. The learning rate was initially 0.001 for the new fully connected layer and was 0.0001 for the rest of the model. To adapt the size change of different data sets, the overall iterations for fine-tuning was set to 300 epochs (9600 iterations) and the attenuation was 0.1 for each 60 epochs.

The advantages of the HRRS pretraining can be summarized as follows,

- 1) Compared to 600 000 iterations of ImageNet pretraining [4], the convergence of the proposed

method is faster than ImageNet pretraining by six times, which greatly reduces pretraining iterations, correspondingly saves training time and computation resources.

- 2) In the experiment, the HRRS based learning strategy is effective to accelerate and improve the convergence of the model which achieve a similar or better performance than ImageNet pretraining. It is valuable to train task-driven model in the remote sensing community.

III. EXPERIMENTS AND ANALYSIS

A. Data Set Description and Experimental Settings

In this section, the proposed scene classification method is evaluated using two publicly available data sets, the UC Merced data set (UCM) [19] and the AID [2]. The UCM was comprised of large optical images extracted from the U.S. Geological Survey, wherein 2100 images with a resolution of 256×256 pixels were selected and manually labeled as belonging to 21 LULC classes. The AID data set is a large-scale data set for aerial scene classification, which contains 10 000 annotated aerial images with a fixed resolution of 600×600 pixels, classed within 30 classes.

The entire learning process contains two main parts: the pretraining phase and the fine-tuning phase. In the pretraining phase, the CNN models included AlexNet [6], VGG-16 [5], VGG-19, and ResNet-18 [4] for comparison, which were trained on ImageNet by 600 000 iterations.

The settings of learning rate and iteration number were introduced in Section II, and the batch size was set to 256 in all experiments. A stochastic gradient descent (SGD) optimizer was used to update the CNN models. The two data sets were randomly split into two groups, 80% for training and 20% for testing. All experiments were performed five times, and the overall results were averaged. Data augmentation was not used in the Experimental section.

Three metrics were adopted to evaluate the convergence properties of the scene classification results: overall accuracy (OA), and the kappa coefficient (kappa).

The proposed method was implemented with an 18-layer structure, of which the depth was equal to ResNet-18. The classification performance of the proposed method and the comparison methods were observed and analyzed. The different learning strategies, including From Scratch, HRRS Pretraining, and ImageNet Pretraining, were compared.

B. Experimental Results and Analysis

The experimental results for the UCM data set and the AID data set are shown in Table II. It can be seen that: 1) ResNet-18 is significantly superior to the early model, AlexNet, by about 2.9% in the UCM and 5.3% in the AID for the OA; 2) The proposed method achieved the best classification with an OA of 98.8% and 94.3% and kappa coefficients of 0.98 and 0.94 using the UCM and AID data sets, respectively, which are higher than all of the comparison methods within the ImageNet pretraining.

For analysis of the effect of learning strategy, a comparison between the accuracy curves of From Scratch, ImageNet

TABLE II
ANALYSIS OF THE EFFECT OF THE PRETRAINING
STRATEGY FOR THE CNN MODELS

Model	UCM		AID	
	OA(%)	Kappa	OA(%)	Kappa
AlexNet	95.47±1.15	0.95±0.01	87.45±0.89	0.87±0.02
VGG-16	96.53±0.96	0.96±0.01	90.25±1.15	0.90±0.02
VGG-19	96.38±1.52	0.96±0.02	89.95±1.53	0.90±0.05
ResNet-18	98.33±0.53	0.98±0.01	92.75±0.34	0.92±0.01
The proposed method with ImageNet pretraining	98.76±0.85	0.98±0.01	93.90±1.15	0.93±0.02
The proposed method with HRRS pretraining	98.81±0.76	0.98±0.01	94.30±0.67	0.94±0.01

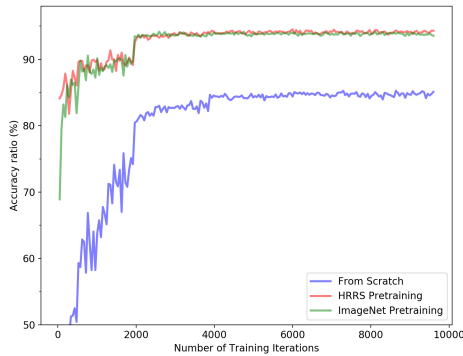


Fig. 3. Comparison of accuracy curves of three kinds of learning strategy on the AID data set: from Scratch, ImageNet Pretraining, and HRRS Pretraining.

Pretraining, and HRRS Pretraining was conducted on the AID data set, which is shown in Fig. 3. As we can see, the method of From Scratch converges slowly and finally reaches a worse convergence, which is lower than the other two results of about 8%. The result of From Scratch is not consistent with the conclusion of the work [20]. The possible reason is that the size of the data set is small, which makes the model overfitting the samples and the generalization ability is relatively poor. By contrast, the two pretraining methods converge fast and reach the accuracy of 90% after 2000 iterations. HRRS Pretraining converges faster and it achieves 80% accuracy in the first iteration. Finally, it reaches a higher convergence degree. The experimental results prove that the pretraining is still important in the remote sensing community, which can accelerate and improve the convergence of the model. Meanwhile, HRRS Pretraining can achieve a 6x speedup and reach a similar or better effect than ImageNet Pretraining, which is valuable to obtain a task-driven model in the field of remote sensing.

IV. CONCLUSION

In this letter, we proposed a deep CNN-based scene classification method with a spatial-scale-aware block and an HRRS imagery-based learning strategy, which is effective for enhancing the classification performance compared to increasing the depth of the models. The spatial-scale-block is not only able to enhance the ability of the spatial representation but also to adaptively recalibrate the channelwise feature responses to

suppress useless feature channels. The HRRS imagery-based learning strategy is presented to obtain an optimized initialization of the model for fine-tuning the model parameters. The experimental results demonstrate that the proposed method is effective for application in HRRS scene classification. In future work, we would apply prior knowledge to improve the performance of the proposed method.

REFERENCES

- [1] G. Cheng, J. Han, and X. Lu, "Remote sensing image scene classification: Benchmark and state of the art," *Proc. IEEE*, vol. 105, no. 10, pp. 1865–1883, Oct. 2017.
- [2] 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.
- [3] F. Hu, G.-S. Xia, J. Hu, and L. Zhang, "Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery," *Remote Sens.*, vol. 7, no. 11, pp. 14680–14707, Nov. 2015.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Las Vegas, NV, USA, Jun. 2016, pp. 770–778.
- [5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. 3rd ICLR*, San Diego, CA, USA, 2015.
- [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. 26th Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [7] Y. Zhong, A. Ma, Y. S. Ong, Z. Zhu, and L. Zhang, "Computational intelligence in optical remote sensing image processing," *Appl. Soft Comput.*, vol. 64, pp. 75–93, Mar. 2018.
- [8] Y. Zhong, X. Han, and L. Zhang, "Multi-class geospatial object detection based on a position-sensitive balancing framework for high spatial resolution remote sensing imagery," *ISPRS J. Photogramm. Remote Sens.*, vol. 138, pp. 281–294, Apr. 2018.
- [9] Y. Yu and F. Liu, "Aerial scene classification via multilevel fusion based on deep convolutional neural networks," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 2, pp. 287–291, Feb. 2018.
- [10] B. Zhao, Y. Zhong, and L. Zhang, "A spectral-structural bag-of-features scene classifier for very high spatial resolution remote sensing imagery," *ISPRS J. Photogramm. Remote Sens.*, vol. 116, pp. 73–85, Jun. 2016.
- [11] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Boston, MA, USA, Jun. 2015, pp. 1–9.
- [12] S. Xie, R. B. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 1492–1500.
- [13] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit. (CVPR)*, Salt Lake City, UT, USA, Jun. 2018, pp. 7132–7141.
- [14] 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. (CVPR)*, Jun. 2009, pp. 248–255.
- [15] W. Zhou, S. Newsam, C. Li, and Z. Shao, "PatternNet: A benchmark dataset for performance evaluation of remote sensing image retrieval," *ISPRS J. Photogramm. Remote Sens.*, vol. 145, pp. 197–209, Nov. 2018.
- [16] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Honolulu, HI, USA, Jul. 2017, pp. 1492–1500.
- [17] H. Zhang, J. Shao, and R. Salakhutdinov, "Deep neural networks with multi-branch architectures are intrinsically less non-convex," in *Proc. AISTATS*, Okinawa, Japan, Apr. 2019, pp. 1099–1109.
- [18] H. Li, C. Tao, Z. Wu, J. Chen, J. Gong, and M. Deng, "RSI-CB: A large scale remote sensing image classification benchmark via crowdsourced data," *CoRR*, vol. abs/1705.10450, 2017.
- [19] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. 18th SIGSPATIAL Int. Conf. Adv. Geographic Inf. Syst.*, Nov. 2010, pp. 270–279.
- [20] K. He, R. Girshick, and P. Dollár, "Rethinking imagenet pre-training," Nov. 2018, *arXiv:1811.08883*. [Online]. Available: <https://arxiv.org/abs/1811.08883>
- [21] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *Proc. Adv. Neural Inf. Process. Syst.*, Montreal, QC, Canada, Dec. 2014, pp. 3320–3328.