# Fusing Local and Global Features for High-Resolution Scene Classification

Xiaoyong Bian, Chen Chen, Long Tian, and Qian Du

*Abstract*—In this paper, a fused global saliency-based multiscale multiresolution multistructure local binary pattern (salM³LBP) feature and local codebookless model (CLM) feature is proposed for high-resolution image scene classification. First, two different but complementary types of descriptors (pixel intensities and differences) are developed to extract global features, characterizing the dominant spatial features in multiple scale, multiple resolution, and multiple structure manner. The micro/macrostructure information and rotation invariance are guaranteed in the global feature extraction process. For dense local feature extraction, CLM is utilized to model local enrichment scale invariant feature transform descriptor and dimension reduction is conducted via joint low-rank learning with support vector machine. Finally, a fused feature representation between salM³LBP and CLM as the scene descriptor to train a kernel-based extreme learning machine for scene classification is presented. The proposed approach is extensively evaluated on three challenging benchmark scene datasets (the 21-class land-use scene, 19-class satellite scene, and a newly available 30-class aerial scene), and the experimental results show that the proposed approach leads to superior classification performance compared with the state-of-the-art classification methods.

*Index Terms*—Codebookless model (CLM), feature representation, image descriptors, rotation invariance, scene classification, saliency detection.

## I. INTRODUCTION

THE recent availability of satellite/aerial images has fostered the development of techniques to classify and interpret high-resolution image scenes with detail structures and different spatial resolutions. In recent years, scene classification has received increasing attention both in academia and real-world application [1]–[5]. The goal of image scene classification is to automatically assign a semantic category to each given image based on some predefined knowledge. Although great effort in image scene classification has been made, it is still a challenging task due to many factors to be considered such as highly complex structures and spatial patterns.

In the past decades, many methods have been presented for image scene classification. The low-level visual feature methods that assume the same type of scene should share certain statistically holistic attributes and have demonstrated their efficiency on classifying image scenes. For example, the scale invariant feature transform (SIFT) [6] was widely used for modeling structural variations in image scenes. In addition, statistical distributions exploitation on certain spatial cues such as color histogram [7], texture information [8] has also been well surveyed. In [9], local structural texture similarity descriptor was applied to image blocks to represent structural texture for aerial image classification. In [10], semantic classification of aerial images based on Gabor and Gist descriptors [11] was evaluated individually. In order to depict the complex scene, the combinations of complementary features are often preferred to achieve improved results. In [12], six different kinds of feature descriptors, i.e., simple radiometric features, Gaussian wavelet features, gray level co-occurrence matrix, Gabor filters, shape features, and SIFT, were combined to form a multiple-feature representation for indexing remote sensing images with different spatial resolutions, and better performance was reported. Recently, local binary pattern (LBP) [13] and completed LBP (CLBP) [14] are also presented. Afterwards, multiscale completed LBP (MS-CLBP) [15] and extended multistructure LBPs (EMSLBP) [16] were adopted for remote sensing image scene classification and competitive results were reported. However, this kind of methods may not able to produce discriminative representation, especially when salient structures in high-resolution image scene often dominate the image category, e.g., the distinct objects such as tennis court, baseball fields, and storage tanks in 21-class land-use scene.

Therefore, more recently, more algorithms for modeling the local variations of structures and evaluating various features via midlevel visual representation have been developed for scene classification. The popular bag-of-visual-words (BoVW) model [17] provides an efficient approach to solve the problem of scene classification, where the image is represented as occurrence histogram of a set of visual words by mapping the local features to a visual vocabulary, where the vocabulary is preestablished after clustering. It should be noted that the original BoVW model ignores spatial and structural information, which may limit its descriptive ability. To avoid this issue, a spatial

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

2

IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING

pyramid matching (SPM) framework was proposed in [18]. This approach partitions an image into fine subregions and computes BoVW histograms of local features in each subregion, and then concatenates the histograms from all subregions to form the SPM representation of an image. However, SPM only considers the absolute spatial arrangement, and the resulting features are sensitive to rotation variations. Later on, several improved land-use scene classification methods with absolute and relative spatial information exploitation are also presented recently. In [19], a spatial co-occurrence kernel, which is general enough to characterize a variety of spatial arrangements, was proposed to capture both the absolute and relative spatial layout of an image. In [20], a multiresolution representation was incorporated into the BoVW model to partition all resolution images into subregions to improve the SPM framework. To achieve rotation invariance, [21] introduces a concentric circle-structured multiscale BoVW model to represent the spatial layout information. In [22], partlet-based method was proposed to achieve efficient VHR image land-use classification. In [23], the improved Fisher vector (IFK) [24] feature coding methods were evaluated for scene classification and reported to achieve better performance. Moreover, unsupervised feature learning [3] from very large local patch features and its variants [25] were explored to obtain better scene classification results.

Nonetheless, all the aforementioned methods still carry little semantic meanings. Currently, deep learning methods achieve impressive results on satellite scene classification [26]–[28]. The two freely available deep convolutional neural networks (CNNs), i.e., OverFeat [29] and CaffeNet [30], may be the most popular for learning visual features for classification. In [31], another architecture, GoogLeNet [32], also showed promising performance for aerial images. In [33], multiscale dense CNNs activations from the last convolutional layer were extracted as local features descriptors and further coded through vector of locally aggregated descriptors [34] and IFK to generate the final image representation. In [35], an efficient stacked discriminative sparse autoencoder (shallow-structured model) is proposed to learn high-level features for land-use classification. In addition, some object detection methods in satellite/aerial scenes via deep learning are also presented [36]–[40]. For all the deep CNN architectures used above, either the global or local features were obtained from the networks' pretrained on image scene datasets and were directly used for classification of image scenes. However, deep CNNs have an intrinsic limitation due to the complicated pretraining process to adjust parameters.

In this paper, we propose a fused feature representation method based on global saliency-based multiscale multiresolution multistructure LBP (salM$^3$LBP) and local codebookless model (CLM), which are utilized to extract global and local features, respectively, to characterize both global structures and local fine details of image scenes. Specifically, the salM$^3$LBP is proposed to extract globally rotation invariant features of image scenes. Then, the local enrichment SIFT (eSIFT) descriptor [41] is employed to extract local features and CLM [42] is chosen to model local features into a discriminative representation. The final representation for an image is achieved by fusing salM$^3$LBP and CLM (salM$^3$LBP–CLM) features. It is noted that CLM describes the patch descriptors by a single Gaussian model,

requiring no pretrained codebook and the subsequent coding. Meanwhile, to alleviate the side effect of background clutter on our approach, we also present a simple yet effective patch sampling method based on saliency detection. Experimental results on three benchmark datasets show that the proposed fused feature representation (salM$^3$LBP–CLM) gains very competitive accuracy compared with state-of-the-art methods.

The main contributions of this paper can be summarized as follows.

1) To the best of our knowledge, it is the first attempt to combine global salM$^3$LBP features and local CLM (eSIFT) to achieve a fused representation for image scene classification. Two different types of features together can effectively mitigate respective shortcomings of global features and local ones.
2) The proposed representation framework is unified in a simple and effective way, which benefits image scene classification.
3) Saliency-based sampling is useful and efficient to exclude background clutter and extracts the representative patches that dominate the image category, which is beneficial to highly cluttered scene (e.g., 21-class land-use scene).

A preliminary version of this work appeared in [16]. This paper extends the earlier work [16] in the following aspects. First, we perform more comprehensive surveys on related works. Second, we proposed the improved LBPs by designing salM$^3$LBP, thereby enhancing the discrimination information among land-use and land-cover (LULC) classes. Third, we develop a fused representation based on global salM$^3$LBP and local CLM with eSIFT descriptor, emphasizing the overall contribution for scene classification, where the weighted coefficient is empirically found by cross-validation strategy. We extensively evaluate our method on three benchmark datasets and comprehensively compare it with the state-of-the-art approaches including deep learning methods, e.g., [27], [28], [31], [33], [49]. Experimental results show that our method outperforms the state-of-the-art methods. Moreover, our approach is flexible to combine with more informative descriptors for further improvement.

The rest of this paper is organized as follows. Section II represents the related works including extended LBP and CLM. Section III describes the proposed salM$^3$LBP–CLM approach in details. Section IV evaluates the proposed approach against various state-of-the-art classification methods on three challenging image scene datasets. Finally, Section V concludes the paper.

## II. RELATED WORK

### A. Extended Local Binary Patterns

LBPs are an effective measure of spatial structure information of local image texture. Given a center pixel and its gray value $x_{0,0}$. Its neighboring pixels are equally spaced on a circle of radius $r$ with the center at location $x_{0,0}$. Suppose the coordinates of the central pixel (CI) are $(0, 0)$ and $p$ neighbors $\{x_{r,n}\}_{n=0}^{p-1}$ are considered. Let $a = 2\pi n/p$, for circular neighborhood, the coordinates of $x_{r,n}$ are $[-r\sin(a), r\cos(a)]$. Then, the LBP is calculated by thresholding the neighbors $\{x_{r,n}\}_{n=0}^{p-1}$ with the center pixel $x_{0,0}$ to generate an $p$-bit binary number; for elliptical neighborhood, let the length of the

minor axis be equal to the radius $r$ of circular neighborhood and set a certain ratio of elliptical major and minor axis as $m$, $X_0 = mr\cos(a + \theta)$, $Y_0 = r\sin(a + \theta)$, then the $x$-coordinate of $x_{r,n}$ is $-X_0 \sin(\theta) + Y_0 \cos(\theta)$, while its $y$-coordinate of $x_{r,n}$ equals $X_0 \cos(\theta) + Y_0 \sin(\theta)$, where four different rotational angles $\theta \in \{0°, 45°, 90°, 135°\}$ in each ellipse are used. Those locations not falling exactly on a pixel are estimated by interpolation. No matter what the neighbor structure, the resulting LBP for $x_{0,0}$ in decimal number can be expressed as follows:

$$\text{LBP}_{p,r} = \sum_{n=0}^{p-1} s(x_{r,n} - x_{0,0})2^n, \ s(x) = \begin{cases} 1, \ x \geq 0 \\ 0, \ x < 0 \end{cases} \quad (1)$$

where the difference $(x_{r,n} - x_{0,0})$ of the neighborhoods against the center pixel characterizes the spatial local structure at center location and is robust to illumination changes because the signs of the differences are utilized. However, the original LBP ignores the magnitudes of the differences. And the rotation-invariant "uniform" LBP (called $\text{LBP}_{p,r}^{\text{riu2}}$) suggested by Ojala *et al.* [13] is defined by

$$\text{LBP}_{p,r}^{\text{riu2}} = \begin{cases} \sum_{n=0}^{p-1} s(x_{r,n} - x_{0,0}), \text{ if } U(\text{LBP}_{p,r}) \leq 2 \\ p+1, \qquad\qquad\qquad\quad \text{otherwise} \end{cases} \quad (2)$$

where

$$U(\text{LBP}_{p,r})$$
$$= |\sum_{n=0}^{p-1} s(x_{r,n} - x_{0,0}) - s(x_{r,\text{ mod }(n+1,p)} - x_{0,0})| \quad (3)$$

where the superscript riu2 denotes the rotation invariant "uniform" patterns that have $U$ values at most 2. Therefore, mapping from $\text{LBP}_{p,r}$ to $\text{LBP}_{p,r}^{\text{riu2}}$ results in only $p+1$ distinct groups of patterns, leading to a much lower histogram representation for the whole image. As the limited abilities of $\text{LBP}_{p,r}$ and $\text{LBP}_{p,r}^{\text{riu2}}$, later on, a CLBP [14] is proposed to use both signs and magnitudes of LBP to extract detailed local structure and texture information for texture classification and claims better performance. More recently, in EM-SLBP [16], the intensity and difference components are both exploited to improve scene classification. The intensity-based descriptors consider the intensity of the CI and those of its neighbors (NI); while for the difference-based descriptors, the radial difference (RD) and angular difference (AD) are computed. Specifically, two rotation-invariant-uniform intensity-based descriptors, CI- $\text{LBP}_{p,r}^{\text{riu2}}$ (abbreviated as CI-$\text{LBP}_{p,r}$) and NI- $\text{LBP}_{p,r}^{\text{riu2}}$ (abbreviated as NI-$\text{LBP}_{p,r}$); two rotation-invariant-uniform difference-based descriptors, RD- $\text{LBP}_{p,r}^{\text{riu2}}$ (abbreviated as RD-$\text{LBP}_{p,r}$) and AD- $\text{LBP}_{p,r}^{\text{riu2}}$ (abbreviated as AD-$\text{LBP}_{p,r}$) are computed, which are different and complementary types of descriptors. More details on the formal definitions of the first three descriptors CI-$\text{LBP}_{p,r}$, NI-$\text{LBP}_{p,r}$, and RD-$\text{LBP}_{p,r}$ can be found in [16]. Obviously, NI-$\text{LBP}_{p,r}$ and $\text{LBP}_{p,r}^{\text{riu2}}$ differ in the selection of thresholding value and NI-$\text{LBP}_{p,r}$ tends to be more discriminative and effective. As for RD-$\text{LBP}_{p,r}$, the objective is to obtain local RD patterns computed with given integer radial displacement $t$ (e.g., $t = 1$), $x_{r,n}$, and $x_{r-t,n}$ according to

the pixel values of pairs of pixels of the same radial direction, therefore, it is more robust to noise. Similarly, AD-$\text{LBP}_{p,r}$ is defined as

$$\text{AD-LBP}_{p,r} =$$
$$\begin{cases} \sum_{n=0}^{p-1} s(x_{r,n} - x_{r,\text{ mod }(n+t,p)}), \text{ if } U(\text{LBP}_{p,r}) \leq 2 \\ p+1, \qquad\qquad\qquad\qquad \text{otherwise} \end{cases}$$
$$s(x) = \begin{cases} 1, \ x \geq \varepsilon \\ 0, \ x < \varepsilon \end{cases} \quad (4)$$

where the AD is computed with given angular displacement $t(2\pi/p)$, where $t$ is an integer such that $1 \leq t \leq p/2$, $x_{r,n}$, and $x_{r,\text{mod}(n+t,p)}$ according to the pixel values of pairs of pixels of $t$ equally spaced pixels on a circular radius $r$, and function $\text{mod}(x, y)$ is the modulus $x$ of $y$. $\varepsilon$ is a threshold value and 1% of the pixel value range.

As stated, the intensity-based features and the difference-based ones are complementary and have powerful discrimination for satellite/aerial scene classification. Let $\text{LBP}_{x,p,r}$ be any of the four local feature descriptors aforementioned, and $\text{LBP}_{x,p,r}(i, j)$ is the extracted LBP pattern of each pixel$(i, j)$, then feature extractor $h_x$ of length $K$ is computed by

$$h_x(k) = \sum_{i=1}^{N} \sum_{j=1}^{M} \delta(\text{LBP}_{x,p,r}(i, j) - k) \quad (5)$$

where $0 \leq k \leq K - 1$, $K = 2^p$ is the number of LBP codes, the subscript $x$ represents circular sampling ("c") or elliptical sampling ("e"), and $\delta(\cdot)$ is the Dirac delta function. $M$ and $N$ are the size of the image.

### B. Codebookless Model

The widely used bag-of-features (BoF) methods such as BoVW and Fisher vector (FV) code an image with a pretrained codebook, where the learned codebook describes the distribution of feature space and makes coding of high-dimensional features possible. However, in the BoVW model, the computational cost scales as the product of the number of visual words and may be less effective. For FV model, the number of Gaussian mixture model (GMM) is often hard to select due to lacking of general criteria as well as high computational cost of matching methods. The most appealing one is CLM [42], which exploits a single Gaussian model to represent image for classification. More specifically, the CLM consists of a Gaussian model and the matching method using Gaussian embedding and the log-Euclidean metric. Due to informative representation with the compact Gaussian model, this representation contains distinct and can be a competitive discriminative information alternative to the BoF methods.

CLM is a direct statistical estimation on sets of dense local features. We first focus on enhancing local features such as SIFT with both raw image information and the relative location and scale of local features within the spatial support of the region, instead of coding local features using large codebooks. Then, a single Gaussian model is employed to simply model the dense local features and further joint low-rank learning with support

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

4           IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING

vector machine (SVM) proposed for high-dimensional feature reduction while respecting the Riemannian geometry structure of Gaussian models, the resultant local features are compact and discriminative. Specifically, three main steps are incorporated in local CLM feature learning. First, dense local features at a dense grid are extracted by eSIFT, please refer to [41] for details on the local descriptor enrichment process.

Second, a single Gaussian model together with two-step metric and two-parameter trade-off is utilized to model eSIFT for discriminative representation. In what follows, we review Gaussian model for image representation in detail. Let $X = \{x_i \in R^{k \times 1}, i = 1, ..., D\}$ be the set of $D$ local features extracted on a dense grid of image. The image can be represented by the following Gaussian model via the maximum likelihood method:

$$\mathcal{N}(x_i|\mu, \Sigma) = \frac{\exp(-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu)}{\sqrt{(2\pi)^k \det(\Sigma)}} \quad (6)$$

where $\Sigma = \frac{1}{D-1} \sum_{i=1}^{D} (x_i - \mu)(x_i - \mu)^T, \mu = \frac{1}{D} \sum_{i=1}^{D} x_i$ are covariance and mean vector matrix, respectively, and $\det(\cdot)$ denotes matrix determinant. Therefore, $\mathcal{N}(\mu, \Sigma)$ is a Gaussian model estimated on a set of local descriptors extracted from input image. Unlike GMMs used in FV, Gaussian models are more efficient to fit the data and informative as well. To fit the data, a two-step metric between Gaussian components is exploited to compute the ground distance. The first step is to embed Gaussian manifold into the space of symmetric positive definite (SPD) matrices. Formally, $\mathcal{N}(\mu, \Sigma)$ is first mapped to an affine matrix A through a continuous function $\pi$, that is,

$$\pi : \mathcal{N}(\mu, \Sigma) \mapsto A = \begin{bmatrix} P & \mu \\ 0^T & 1 \end{bmatrix} \quad (7)$$

in the affine group, $A_k^+ = \{(\mu, P)|\mu \in R^{k \times 1}, P \in R^{k \times k}, \det(P) > 0\}$ is an $k$-dimensional element and $\Sigma = PP^T$ is the Cholesky factorization of $\Sigma$. A is mapped to an SPD matrix S through the successive function $\gamma: A \mapsto S = AA^T$. Consequently, $\mathcal{N}(\mu, \Sigma)$ is uniquely designed as a $(k+1) \times (k+1)$ SPD matrix, namely,

$$\mathcal{N}(\mu, \Sigma) \sim S = \begin{bmatrix} \Sigma + \mu\mu^T & \mu \\ \mu^T & 1 \end{bmatrix} \quad (8)$$

where S is an embedding matrix and the space of $(k+1) \times (k+1)$ SPD matrices $S_{k+1}^+$ is a Lie group that forms a Riemannian manifold. Then map $S_{k+1}^+$ into its corresponding Lie algebra $S_{k+1}$, the vector space of $(k+1) \times (k+1)$ symmetric matrices, a linear space by using the Log-Euclidean metric. Moreover, on one hand, as observed in [42], mean vector and covariance matrix in the embedding matrix may have different effects, meanwhile, their dimensions and order of magnitude of each dimension may vary. Thus, a parameter $\beta\,(\beta > 0)$ is introduced to balance the effect between mean and covariance of Gaussian in (8); on the other hand, with consideration of the observations that the maximum likelihood estimator of covariance is susceptible to noise interference and ought to be improvable by eigenvalue shrinkage [43], hence another parameter $\rho$ is introduced in the normalization of covariance matrix, where

eigenvalues power normalization (EPN) is applied to estimate covariance matrices. Therefore, (7) and (8) can be rewritten, respectively, as

$$\pi(\beta) : \mathcal{N}(\mu, \Sigma) \mapsto A = \begin{bmatrix} P & \beta\mu \\ 0^T & 1 \end{bmatrix} \quad (9)$$

$$\mathcal{N}(\mu, \Sigma) \sim S(\beta, \rho) = \begin{bmatrix} \Sigma^\rho + \beta^2 \mu\mu^T & \beta\mu \\ \beta\mu^T & 1 \end{bmatrix}. \quad (10)$$

It is proved that the embedding matrix $S(\beta, \rho)$ is still positive definite as $\Sigma^\rho$ being an SPD matrix. Note that EPN is characterized for robust estimation of covariance matrices in Gaussian setting for the case of high-dimensional features and comparison of Gaussians in Gaussian embedding and the Log-Euclidean metric. The matrix $S(\beta, \rho)$ can be further embedded into a linear space by matrix logarithm:

$$G(\beta, \rho) = \log(S(\beta, \rho)). \quad (11)$$

So an SPD matrix S is one-to-one mapped to a symmetric matrices G (for simplicity, omitting the parameters $\beta, \rho$) which lies in a linear space, and the geodesic distance between two Gaussian models $\mathcal{N}_i = \mathcal{N}(\mu_i, \Sigma_i)$ and $\mathcal{N}_j = \mathcal{N}(\mu_j, \Sigma_j)$ is

$$\text{dist}_{\mathcal{N}_i, \mathcal{N}_j} = ||G_i - G_j||_F \quad (12)$$

where $F$ is the Frobenius norm and distance (12) is known as decoupled so that $G_i$ and $G_j$ can be computed separately and adopted in a linear classifier.

Finally, joint low-rank learning with SVM strategy is performed to reduce high dimensional $(> 10^4)$ local eSIFT features in this model in order to achieve compact CLM. By introducing a low-rank transformation matrix $L \in R^{d \times r}, r << k^2, d = (k+1) \times (k+2)/2$, the geodesic distance (12) can be modified as

$$\text{dist}_{\mathcal{N}_i, \mathcal{N}_j} = ||L^T(f_i - f_j)||_2, \; s.t. \; L^T L = I \quad (13)$$

where $f_i$ and $f_j$ are the unfolding vectors of two Gaussian models $\mathcal{N}_i$ and $\mathcal{N}_j$, respectively. Motivated by joint optimization of dimensionality reduction with classifier, here the low-rank learning is jointly optimized with a linear SVM in an SPM framework. Readers can refer to [42] for more details. As a result, the reduced local features ought to be more informative while reducing computational cost for classification hereafter. Note that CLM differs in the way other ad-hoc local feature methods are handled and is discriminative.

Some other feature representation methods for scene classification can be found in [3], [4], [7], [15]–[17], [20], [21], [25], [27], [28], [31], [33], [42], and [44], whereas our work is most related to [42] and [44]. In [42], Wang *et al*. evaluated the CLM features with various local descriptors in image databases, showing impressive classification performance, whereas we fuse CLM with global descriptor for scene classification. In [44], local and global features are employed and fused based on minimal residuals for classification; however, sparse representation may be unstable to some extent. With respect to deep learning based methods, in Castelluccio *et al*. [31] and Hu *et al*. [33] employed direct pretrained CNNs or fine-tuned pretrained CNNs for scene

classification and state-of-the-art results were reported. In contrast to non-CNNs, CNNs have to tune so many parameters of all layers at the same time which is often time-consuming.

## III. PROPOSED FEATURE REPRESENTATION FRAMEWORK

Driven by the success of extended LBP and codeless model (CLM) in computer vision communities, we propose an effective and efficient image presentation approach for high-resolution image scene classification, that is, the fusion of salM$^3$LBP and CLM. The salM$^3$LBP is used as global feature descriptor, while patch-based eSIFT as local feature descriptor and modeled by CLM. Then the salM$^3$LBP and CLM (eSIFT) are fused as scene descriptor for classification, represented as salM$^3$LBP–CLM. Here, we describe our scene classification framework. The overall framework of the proposed salM$^3$LBP–CLM is shown in Fig. 1. As depicted in Fig. 1, the method consists of four parts.

1) First, converting the images from RGB color space to YCbCr color space to obtain a grayscale image using the Y component and multiple down-sampled scales of the original image are obtained. Meanwhile, a saliency-based patch sampling is proposed to enhance our salM$^3$LBP when images are heavily cluttered. Note that saliency detection is just for some classes in the cluttered scenes.

2) Each scale image of the dataset is fed to extract global features, where multiple resolution and two types of structure (circular/elliptical sampling) are applied via the proposed CI-LBP$_{p,r}$, NI-LBP$_{p,r}$, RD-LBP$_{p,r}$, and AD-LBP$_{p,r}$ descriptors. Note that the extracted micro/macrostructure features (salM$^3$LBP) are rotation invariant in circular sampling and rotation invariance of those from elliptical sampling can be derived by averaging the histograms over different rotational angles.

3) Then, dense local features are extracted by local eSIFT and further modeled by CLM in Gaussian setting and joint low-rank learning with SVM manner. Moreover, the reduced local features CLM and above global features (salM$^3$LBP) are fused to generate global image representation, and output the distribution of these features as the scene descriptor.

4) Finally, kernel-based extreme learning machine (KELM) [45] is adopted for scene classification and label assignment. Moreover, accuracy evaluation is extensively conducted.

### A. Saliency Detection

We first present a simple yet effective saliency patch sampling method based on unsupervised saliency detection [46]. This method focuses attention on the image regions that are most informative and dominate the image category, to represent the scene information in the image and can detect the visual saliency patches in the global and local perspectives. That is, this method unifies global and local saliencies by measuring the dissimilarities among image patches according to the "repetition suppression principle" in the area of brain cognition. Specifically, our saliency detection method consists of two steps: novel items search and partial background removal. In the first step we localize in novel items in cluttered scenes

and then determine the bounding box surrounding the novel items (foreground). Next, we adaptively expand bounding box to accommodate some background regions based on size and intensity variance of the area inside the bounding box. Specifically, we first choose the pixels with consistency score higher than object threshold in the salient map as salient region and fit a bounding box in the salient region, then heuristically expand it with a centered larger size of width and height such as [1.2 1.5] until image boundary and smaller intensity variance ratio such as 0.8, which is a subimage with less clutter and helps to guide the sampling. That is to say, some nonsalient regions according to the scenes are sampled at the same time, for on one hand, not all scene images satisfy the assumption that the salient regions usually correspond to the scene; on the other hand, neighboring regions of object can serve as the context and may be helpful for classification. Then, the area outside bounding-box is removed for classification. The purpose here is to remove the interference of background for local and global feature extraction. Based on aforementioned observations, we define novel items as coming from those regions that hold dissimilarities with both global and local properties, and saliency of a patch $I_i$ is computed by

$$\text{sal}(I_i) = \sum_f \sum_{k=1}^{N^f} H(W_{f,k}) \cdot \varphi_{f,k}, \ f \in \{\text{color, texture}\} \tag{14}$$

$$W = \{\{w_k^{\text{color}}\}_{k=1}^{N^f}; \{w_k^{\text{texture}}\}_{k=1}^{N^f}\} \tag{15}$$

where $H(\cdot)$ means the histogram of visual words and $\varphi_{f,k}$ a weighed factor for each visual word. $w_k^{\text{color}}, w_k^{\text{texture}}$ denote the $k$th color and texture words, respectively. $N^f$ is the number of the quantized words in the image. We compute the saliency regions of some cluttered scenes in the considered dataset. Fig. 2 shows the saliency detection results for different LULC classes in the 21-class, 19-class, and 30-class scenes. In this paper, we first adopt saliency detection to the three scene datasets (21-class, 19-class, and 30-class) with some heavy background clutters, for instance, residential, tenniscourt, bridge, river, playground, and stadium, as they are simultaneous visually dominating objects in the scenes, these classes benefit most from the proposed saliency sampling. It should be noted that we do not perform saliency detection on images with less background clutter and image scene where both foreground and background are valuable for scene representation. For example, most of natural scene types are easy to be distinguished, thus these types are saliency-free, whereas for those man-made scene types such as residential, footballField, and storagetanks, saliency detection can be done before feature extraction. Furthermore, we also use multiscale saliency to measure the saliency of a patch in a multiscale image. The proposed saliency detection is also applied in the subfeatures.

### B. Global and Local Feature Extraction

For global feature extraction, the proposed salM$^3$LBP approach is applied to the scene (or saliency-based scene) to extract extended LBP feature histogram according to a set of parameters such as sampling point $p$, radius $r$, the ratio of elliptical major and minor axis $m$, rotational angle $\theta$, etc. In

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

6                                                              IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING
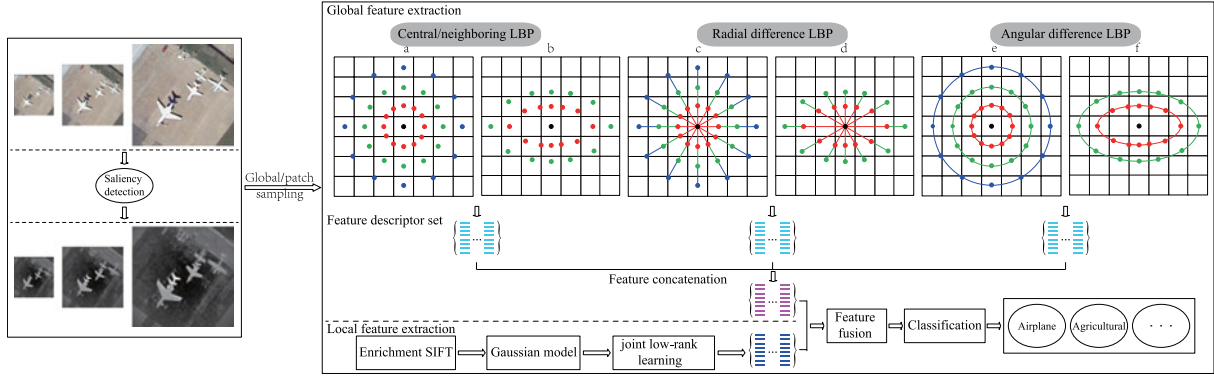


Fig. 1.    Overall architecture for the proposed approach.



Fig. 2.    Saliency detection results for three different scene images. From top to bottom, the first two lines: 21-class, the middle two lines: 19-class and the last two lines: 30-class.

order to incorporate more information for image description, the multiresolution and multistructure sampling is enforced to the global descriptor. Specifically, two types of LBP sampling (circular and elliptical structures) with multiple radii and fixed sampling points are considered to extract complementary features. Formally, for an image patch, three-coupled global descriptors (i.e., circular/elliptical NI-LBP, circular/elliptical RD-LBP, and circular/elliptical AD-LBP) can be employed to extract global features and result in multiple feature vectors that can capture texture, structure information, and spatial patterns such as flat area, spot, corner, edge. Meanwhile, multiple scales are also considered in salM$^3$LBP; thus, we alter image scales in different levels to capture both microstructure and macrostructure properties. More precisely, multiple scale ($l$), multiple radii

resolution ($r$), and multiple structure ($t$) are used in global feature extraction process, then the global preliminary features with size of $c \times l \times r \times t$ ($c$ is a constant) are generated. Note that the extracted features from circular sampling are isotropic and rotation invariant; however, those extracted from elliptical sampling are anisotropic and should be averaged over different rotational angles to derive rotation invariance. The reason is that average anisotropic histogram is insensitive to local image fluctuation such as rotation and its use as statistical feature of each image is globally invariant to these changes. It is mentionable that, as observed in [47], the proportions of the uniform patterns of AD-LBP are too small to provide a reliable and meaningful description of texture. Consequently, we mainly focus on other global descriptors instead of AD-LBP. The global features are directly stacked as a final feature vector.

For local feature extraction, we use local eSIFT descriptor to obtain high-dimensional local features from the scene (or saliency-based scene) aforementioned, followed by CLM to produce reduced local feature vector. Specifically, Given an image $I$, $B$ blocks in SPM framework, $D$ dense local eSIFT features can be extracted, which is fed to CLM to fit a single Gaussian model that will be embedded into a vector space, then its SPM representation is obtained such that a joint low-rank learning with SVM is allowed to perform, as a result, a set of reduced discriminative local features are achieved, shown in Fig. 3.

As illustrated in Fig. 3, local CLM (eSIFT) feature is discriminatively represented by matrix multiplication between a learned Gaussian model and a low-rank transformation matrix. Compared with histogram and covariance, Gaussian model is more informative and fitting of Gaussian models does not bring high computational cost, unlike GMMs used in FV.

### C. Classification by Fusing Features From Two Scenarios

As described above, the proposed scenario (I): salM$^3$LBP features are obtained globally (in the whole image) and scenario (II): CLM (eSIFT) features, for notational simplicity, called CLM, are achieved locally (extracted from overlapping local image patches). Each type of features (i.e., global salM$^3$LBP and local CLM) reflects various properties and has its own meanings such as LBP feature reveals the local image texture (e.g., edges, corners, shapes, etc.), and CLM feature captures variations, relations, etc. Based on the above-mentioned characteristics, we

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

BIAN *et al.*: FUSING LOCAL AND GLOBAL FEATURES FOR HIGH-RESOLUTION SCENE CLASSIFICATION 7
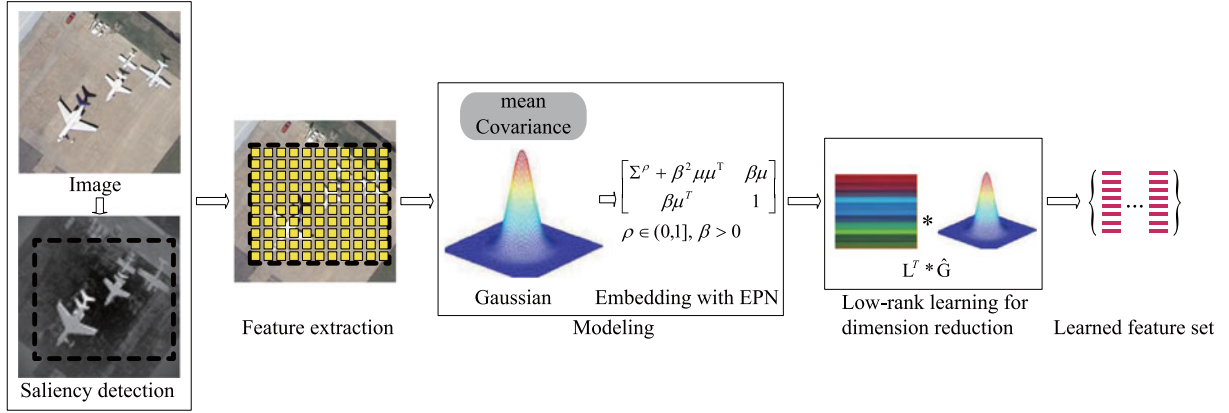
Fig. 3. CLM representation.

propose to fuse two scenarios, i.e., the global and local features together, to characterize both global structures of scenes and local fine details for scene classification. Cross-validation strategy is employed for tuning the optimal parameter (i.e., $\eta$) using training data. In the proposed global-local fusion strategy, each test image $y$ is generated two types of features, i.e., $F_G$ and $F_L$ representing global and local features, respectively, and feature-level fusion is formulated as

$$F_f(y) = \eta F_G(y) + (1 - \eta) F_L(y) \qquad (16)$$

where the feature $F_f$ is a final feature representation of test image $y$. Class label assignment on the test image is conducted using standard distance metrics such as KELM. As illustrated in Fig. 1, our fused scene descriptor promises a substantial performance enhancement compared to individual decisions.

## IV. EXPERIMENTS

Three challenging remote sensing scene datasets are utilized for extensive evaluations of the proposed approach in the experiments: 21-class land-use scene, 19-class satellite scene, and 30-class aerial scene. In the experiments, the KELM is adopted for classification. The classification performance of the proposed method is compared with the state-of-the-arts in the literature.

### A. Image Data and Experimental Setup

The first dataset is the well-known UC-Merced land-use dataset consisting of 21 land-use classes and each class contains 100 images of size $256 \times 256$ pixels with spatial resolution of 30 cm (1 foot). It is up to now the first public land-use scene image dataset with ground truth, which is manually extracted from aerial orthoimagery downloaded from the United States Geological Survey (USGS) National Map. This is a challenging dataset due to a variety of spatial patterns and some highly overlapped classes, e.g., dense residential, medium residential, and sparse residential that mainly differ in the density of structures, which make the dataset difficult for classification. This benchmark dataset has a large geographical scale. For more information, see [1] and visit http://vision.ucmerced.edu/datasets.

The second dataset used in our experiments is composed of a 19-class satellite scene dataset. It consists of 19 classes of high-resolution satellite scenes. Each class has 50 images,

TABLE I
SCENE CLASSES AND THE NUMBER OF IMAGES PER CLASS IN 30-CLASS AERIAL SCENE DATASET

| Name | #images | Name | #images | Name | #images |
|---|---|---|---|---|---|
| Airport | 360 | farmland | 370 | port | 380 |
| bare land | 310 | forest | 250 | railway station | 260 |
| baseball field | 220 | industrial | 390 | resort | 290 |
| Beach | 400 | meadow | 280 | river | 410 |
| bridge | 360 | medium residential | 290 | school | 300 |
| center | 260 | mountain | 340 | sparse residential | 300 |
| church | 240 | park | 350 | square | 330 |
| commercial | 350 | parking | 390 | stadium | 290 |
| dense residential | 410 | playground | 370 | storage tanks | 360 |
| desert | 300 | pond | 420 | viaduct | 420 |

with sizes of $600 \times 600$ pixels. This data set is a challenging one because all these scenes are extracted from very large satellite images on Google Earth, where the illumination, appearances of objects, and their locations vary significantly, with frequent occlusions. For more information, see [18] and visit http://dsp.whu.edu.cn/cn/staff/yw/HRSscene.html.

The third dataset is made up of 30 aerial scene types and all the images are collected from the Google Earth. The number of different aerial scene types varies a lot, see Table I, from 220 up to 420, and a total number of 10 000 images. As far the 30-class scene is the largest annotated aerial image datasets which is available online from http://www.lmars.whu.edu.cn/xia/AID-project.html. Since the images in this scene are from different remote imaging sensors and extracted at different time and seasons under different imaging conditions, which increases the intraclass diversities of the dataset as well as low interclass dissimilarity, thus this brings more challenges for scene classification than the single source images. Different from two previous image scenes, the 30-class has multiresolutions ranging from about 8 m to about half a meter, and thus the size of each aerial image is fixed to be $600 \times 600$ pixels to cover a scene with various resolutions.

Note that the original images in these three datasets are color images; the images are converted from the RGB color space to the YCbCr color space, and the Y component (luminance) is used for scene classification.

The parameter settings in our experiments are given as follows.

1) For training set generation, we adopt two different settings for each tested dataset in the supervised classification process. For 21-class dataset, the ratios are set to be 50% and 80%. For 19-class dataset, the ratios are fixed at 40% and 60%, while for 30-class dataset, we fix the ratio of the number of training set to be 20% and 50%, and the left for testing. We randomly split the datasets into training sets and testing sets for evaluation.

2) For features inclusion and transformation, different local feature descriptors (i.e., global salM$^3$LBP and local CLM) are employed to exploit the discriminative feature information and fuse for scene classification. Specifically, in global LBP feature extraction, for multiresolution analysis, the sampling points are empirically fixed as $p = 16$ for all three scene datasets, and 8 radii (i.e., $r = [1 : 8]$) are used for the salM$^3$LBP; for multiscale analysis, three scales of down-sampled images including $\{1, 1/2, 1/3\}$ with 1 being the original image are considered. Meanwhile, both circular and elliptical structures are exploited for extended LBP sampling, while for elliptical sampling, the LBP histograms from four different rotational angles $\theta \in \{0°, 45°, 90°, 135°\}$ in each ellipse are extracted, then averaged and stacked with isotropic (circular-sampling) features as an LBP feature vector. All this results in a relatively small feature size of $2 \times 2 \times 3 \times 8 \times 18$ (1728) global features. As for local CLM feature extraction, dense local descriptor eSIFT using a single patch size is employed to generate high-dimensional local features that further can be reduced by joint low-ranking learning with linear SVM. Consequently, a final feature representation via fusing global salM$^3$LBP and local CLM is achieved. The local patch size of 32, 32, and 64 on a dense grid of step size 2 is utilized for 21-class, 19-class, and 30-class scenes, respectively, and typically found that achieves the respective best classification performances.

3) For classification, we report the overall accuracy (OA), kappa statistic ($\kappa$), standard deviation (SD), confusion matrix, and computational time (in seconds). OA is defined as the number of correctly predicted images divided by the total number of predicted images. It is an effective measure to reveal the classification performance on the whole dataset. Confusion matrix is a specific table layout that allows direct visualization of the performance on each class. Each column of the matrix represents the instances in a predicted class, and each row represents the instances in an actual class. To compute the OA and $\kappa$, ten individual runs are conducted and the average results are reported as the means and SDs of OA and $\kappa$. To compute the confusion matrix, we fix the training set by choosing the same images for fair comparison on each datasets and the ratio of the number of training set of the 21-class dataset, the 19-class dataset, and 30-class dataset to be 50%, 40%, and 20%, respectively, whereas for plotted classification performance, 80%, 60%, and 50% are used, respectively.

4) For performance comparison, some strongly related low-level feature methods including BoVW, SPM, multiscale

**TABLE II**
**OVERALL ACCURACY (%) AND STANDARD DEVIATION FOR THE DIFFERENT METHODS WITH DIFFERENT TRAINING RATIOS ON THE 21-CLASS LAND-USE DATASET**

| Methods | 80% labeled samples per class | 50% labeled samples per class |
|---|---|---|
| Proposed salM$^3$LBP–CLM | $95.75 \pm 0.80$ | $\mathbf{94.21} \pm 0.75$ |
| Proposed salM$^3$LBP | $93.14 \pm 1.00$ | $89.97 \pm 0.85$ |
| Proposed M$^3$LBP | $90.95 \pm 1.03$ | $87.49 \pm 1.41$ |
| salCLM (eSIFT) | $94.52 \pm 0.79$ | $92.93 \pm 0.92$ |
| CLM (eSIFT) | $93.62 \pm 0.85$ | $91.88 \pm 1.06$ |
| Combing Scenarios I and II [33] | $\mathbf{98.49}$ | |
| Fine-tuning GoogleNet [31] | $97.10$ | |
| CaffeNet [49] | $95.02 \pm 0.81$ | $93.98 \pm 0.67$ |
| GoogLeNet [49] | $94.31 \pm 0.89$ | $92.70 \pm 0.60$ |
| VGG-VD-16 [49] | $95.21 \pm 1.20$ | $94.14 \pm 0.69$ |
| OverFeat [28] | $90.91 \pm 1.19$ | |
| MS-CLBP+FV [48] | $93.00 \pm 1.20$ | $88.76 \pm 0.79$ |
| Gradient boosting random CNNs [27] | $94.53$ | |
| Partlets-based [22] | $91.33 \pm 1.11$ | |
| Multifeature concatenation [50] | $92.38 \pm 0.62$ | |
| Pyramid of spatial relations [51] | $89.10$ | |
| Saliency-guided feature learning [25] | $82.72 \pm 1.18$ | |
| Unsupervised feature learning [3] | $81.67 \pm 1.23$ | |
| BoVW [31] | $76.81$ | |

completed LBP with FV representation (MS-CLBP+FV) [48], multifeature concatenation and deep CNNs [31], [33], [49] have been implemented. The reader is referred to those papers for additional details. As the basic classifier, KELM is adopted for classification due to its general superiority and low computational cost.

5) For implementation details, to make the comparisons as meaningful as possible, we use the same experimental settings in [1], [18], [31], [33], and [49] for 21-class, 19-class, 30-class and all results are originally reported. It should be noted that each sample is normalized to be zero mean, unit SD, and all the results are reported over ten random partitions of the training and testing sets. All the implementations were carried out using MATLAB R2016a in a desktop PC equipped with an Intel Core i7 CPU (3.4 GHz) and 32 GB of RAM.

## B. Classification of 21-Class Land-Use Scene

We perform a comparative evaluation of the proposed salM$^3$LBP–CLM approach against several state-of-the-art scene classification methods mentioned above on 21-class land-use scene, as shown in Table II. As can be seen from Table II, our method consistently outperforms almost all other scene classification approaches except for two recent CNNs [31], [33] based on the average accuracies and SD obtained by ten trials of random partition of this dataset, with an increase in OA of 0.54%, 0.07% over the second best method, VGG-VD-16, using 80%, 50% labeled samples per class as training ratio on 21-class test set, respectively. This is due to the fact that proposed global salM$^3$LBP feature and local CLM are complementary and fused together to contributing to better performance. Moreover, local CLM features are modeled by a single Gaussian model while preserve the geometry structure of Gaussian models. To our best
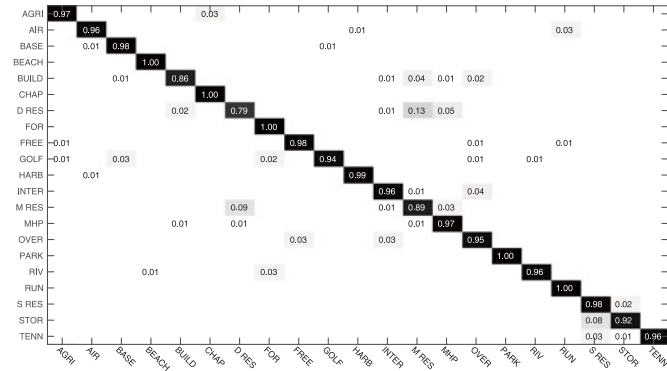
This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

BIAN *et al.*: FUSING LOCAL AND GLOBAL FEATURES FOR HIGH-RESOLUTION SCENE CLASSIFICATION

9

Fig. 4. Confusion matrix obtained by the proposed salM$^3$LBP–CLM on 21-class land-use scene.
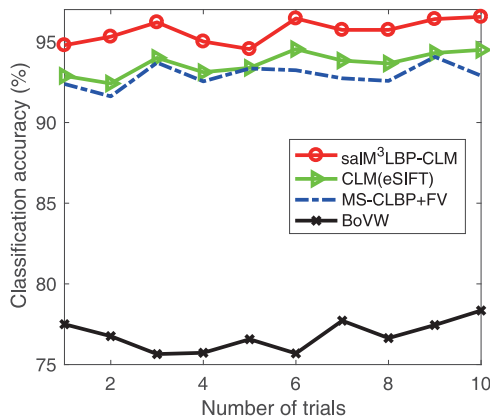


Fig. 5. Classification performance on 21-class land-use scene.

knowledge, this classification result is remarkable on this data set, which adequately shows the effectiveness and superiority of the proposed approach for land-use scene classification. Fig. 4 shows the confusion matrix for the proposed salM$^3$LBP–CLM on 21-class test set. From the confusion matrix, we can observe the consistent phenomenon that our approach obtains a clean confusion matrix. Most of scene types on this scene dataset can achieve the classification accuracy close to or even equal to 1 using the proposed approach, some of them are natural scene types and thus easy to be differentiated, for instance, beach, chaparral, forest, etc. Some difficult classes such as airplane, buildings, dense residential, medium residential, sparse residential, and tennis court are all improved by the proposed salM$^3$LBP–CLM. Meanwhile, relatively high accuracies for some compound objects with spatial recurrent patterns such as harbor, golf court, and storage tanks are also achieved. In addition, the subfeatures of salM$^3$LBP and salCLM (eSIFT) all obtain relatively good accuracies, which validates the effectiveness of our saliency-based subfeature methods.

Furthermore, we plot the classification performance of BoVW, MS-CLBP+FV, and the proposed salM$^3$LBP–CLM corresponding to ten different trials with randomly selected training samples as shown in Fig. 5. It is apparent that the proposed salM$^3$LBP–CLM consistently outperforms the subfeature CLM (eSIFT), MS-CLBP+FV, BoVW, which verifies the

TABLE III
OVERALL ACCURACY (%) AND STANDARD DEVIATION FOR THE DIFFERENT METHODS WITH DIFFERENT TRAINING RATIOS ON THE 19-CLASS SATELLITE DATASET

| Methods | 60% labeled samples per class | 40% labeled samples per class |
|---|---|---|
| Proposed salM$^3$LBP–CLM | 96.38 ± 0.82 | 95.35 ± 0.76 |
| Proposed salM$^3$LBP | 92.58 ± 0.89 | 89.74 ± 1.84 |
| Proposed M$^3$LBP | 91.65 ± 1.14 | 85.47 ± 0.95 |
| salCLM (eSIFT) | 95.92 ± 0.95 | 93.81 ± 0.91 |
| CLM (eSIFT) | 94.82 ± 1.03 | 92.54 ± 1.02 |
| Combing Scenarios I and II [33] | **98.89** | |
| CaffeNet [49] | 96.24 ± 0.56 | 95.11 ± 1.20 |
| GoogLeNet [49] | 94.71 ± 1.33 | 93.12 ± 0.82 |
| VGG-VD-16 [49] | 96.05 ± 0.91 | **95.44** ± 0.60 |
| Multifeature concatenation [51] | 94.53 ± 1.01 | |
| MS-CLBP+FV [48] | 94.32 ± 1.02 | |
| MS-CLBP+BoVW [48] | 89.29 ± 1.30 | |
| Bag of SIFT [51] | 85.52 ± 1.23 | |
| SIFT+LTP-HF+color histogram [7] | 93.6 | |

advantages of salM$^3$LBP–CLM as compared to its counterparts. In addition, we observe that MS-CLBP+FV is much better than BoVW because FV contains more information than a single histogram representation in BoVW. We empirically found that the proposed salM$^3$LBP–CLM achieves the highest accuracy at a small weight of 0.1 for global salM$^3$LBP feature and 0.9 for local CLM, which indicates local feature is assigned with a larger weight, leading to better discrimination. On the other hand, when compared with CNNs, it can be found that our methods are comparable to or slightly lower than two recent CNNs [31], [33], most likely in part due to the CNNs for learning finer grained discriminative features for classification. However, with the increasing depth of the network to the multilayer CNNs architecture, the classification accuracy may oscillate slightly because the deep network has more parameters to train, and the limited number of training samples restricted the performance of the deep network. Thus, the proposed salM$^3$LBP–CLM is very competitive in terms of the classification performance and speed (see Table V) trade-off.

### C. Classification of 19-Class Satellite Scene

In order to further measure the scene classification performance of the proposed approach, we compare the classification accuracies with several state-of-the-art classification methods including deep network [33], [48], [49], [51], [7] in 19-class satellite scene. As in 21-class scene experiments, our salM$^3$LBP–CLM yields highly comparable accuracies with OA of 96.38%, 95.35% using 60%, 40% labeled samples per class, respectively. As can be seen from Table III, the proposed salM$^3$LBP–CLM is consistently better than almost all the others; however, the supremacy of the salM$^3$LBP–CLM is challenged by [33] who have demonstrated that using the combined CNNs in two scenarios can produce better classification performance, this is because deep CNN features are more discriminative; however, the CNN feature learning is computationally expensive. Moreover, there is an exception that our salM$^3$LBP–CLM

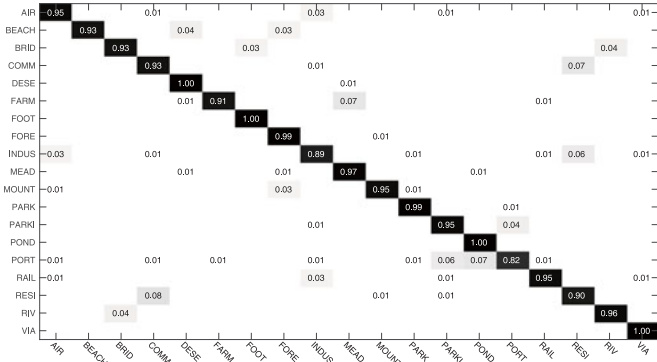This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

10      IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING



Fig. 6. Confusion matrix obtained by the proposed salM$^3$LBP–CLM on the 19-class satellite scene.

TABLE IV
OVERALL ACCURACY (%) AND STANDARD DEVIATION FOR THE DIFFERENT METHODS WITH DIFFERENT TRAINING RATIOS ON THE 30-CLASS AERIAL DATASET

| Methods | 50% labeled samples per class | 20% labeled samples per class |
|---|---|---|
| Proposed salM$^3$LBP–CLM | **89.76** ± 0.45 | **86.92** ± 0.35 |
| Proposed salM$^3$LBP | 87.59 ± 0.38 | 82.31 ± 0.19 |
| Proposed M$^3$LBP | 84.80 ± 0.56 | 80.69 ± 0.33 |
| salCLM (eSIFT) | 88.41 ± 0.63 | 85.58 ± 0.83 |
| CLM (eSIFT) | 87.33 ± 0.68 | 84.21 ± 0.89 |
| CaffeNet [49] | 89.53 ± 0.31 | 86.86 ± 0.47 |
| GoogLeNet [49] | 86.39 ± 0.55 | 83.44 ± 0.40 |
| VGG-VD-16 [49] | 89.64 ± 0.36 | 86.59 ± 0.29 |
| MS-CLBP+FV | 86.48 ± 0.27 | |
| BoVW | 78.66 ± 0.52 | |



Fig. 7. Classification performance on 19-class satellite scene.



Fig. 8. Confusion matrix obtained by the proposed salM$^3$LBP–CLM on the 30-class aerial scene.

is slightly lower (0.09%) than the best method VGG-VD-16 when 40% training ratio is available. For this case, the pre-trained CNN may be more effective with limited training data. In general, the higher the feature level, the better the performance, which indicates mid-/high-level feature methods may be superior to low-level ones. Fig. 6 shows the confusion matrix for the proposed salM$^3$LBP–CLM. From the confusion matrix on 19-class satellite scene, we can see that most of classes are easily distinguished from others that the classification accuracies of almost all classes are above 0.92 by our salM$^3$LBP–CLM. The major confusion occurs between class 9 (i.e., industrial) and class 17 (i.e., residential), or class 13 (i.e., parking) and class 15 (i.e., port), for their similar structures and spatial patterns. Our analysis also shows that the subfeatures salM$^3$LBP and salCLM (eSIFT) can also give relatively good results, as is expected.

Fig. 7 further illustrates the results with ten different trials with random training samples, where the proposed salM$^3$LBP–CLM gains the highest accuracy and local feature can present better than global feature in this satellite scene dataset because the subfeature CLM still dominates the representation in the presence of the large scale of satellite images. As a consequence, a relatively larger weight of 0.15 for global salM$^3$LBP feature and 0.85 for local CLM is found by cross-validation strategy using training data. Therefore, our method has improved performance for high-resolution scene classification.

## D. Classification of 30-Class Aerial Scene

The proposed approach is also carried out on 30-class aerial scene with different resolutions, high intraclass diversity, and low interclass dissimilarity. The classification results of the proposed methods and baseline algorithms for 30-class aerial scene with different training ratios are summarized in Table IV. We observe that the salM$^3$LBP–CLM is superior to the others by a medium margin and consistently achieves improvements beyond the state-of-the-art (e.g., CNN-based methods) with relatively large amount of training data. This result is possibly explained by the fact that the salM$^3$LBP–CLM has the ability to learn discriminative features. Specifically, our salM$^3$LBP–CLM gains a slightly better margin of OA improvements with 0.12%, 0.06% over the second best algorithm, VGG-VD-16 and CaffeNet using 50% and 20% training ratios, respectively. This holds for interpreting the consistency of the proposed salM$^3$LBP–CLM. Although our approach performs better than the others, the proposed method needs more informative features to further enhance the representation power. Fig. 8 shows the confusion matrix for the proposed salM$^3$LBP–CLM on 30-class test set. As in the previous confusion matrices, we can see that most scene types can achieve the classification accuracy more than 0.9 by the proposed method, some of them are natural scene types and thus easy to be partitioned, even the accuracies of these classes such as dense residential, medium

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

BIAN *et al.*: FUSING LOCAL AND GLOBAL FEATURES FOR HIGH-RESOLUTION SCENE CLASSIFICATION
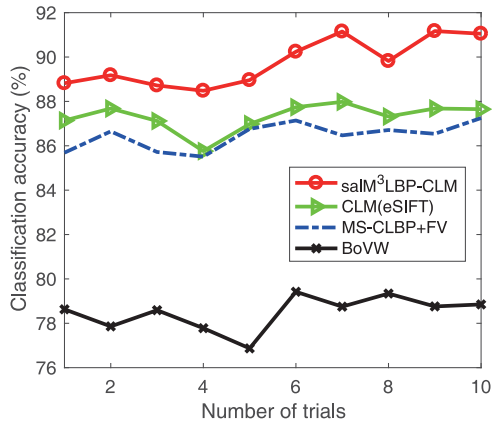
11

Fig. 9.  Classification performance on 30-class aerial scene.

residential, sparse residential are increased to some extent. The most difficult scene types are almost newly added types, for instance, school (0.54), resort (0.57), square (0.63), and center (0.73). It is interesting from Fig. 8 that two of above four classes are improved by our method except for class resort (0.6) and square (0.63) compared with the accuracies of the same classes from the confusion matrix of [49], which demonstrates that the proposed salM$^3$LBP–CLM is effective. The obvious confusion is between resort and park, dense residential and school, which have the similar appearances and may contain the same image clutter such as green belts and buildings, respectively, and thus are easily confused. Experimental results with ten trails with random training samples are illustrated in Fig. 9. It is apparent that the proposed salM$^3$LBP–CLM outperforms the subfeature CLM (eSIFT), MS-CLBP+FV, and BoVW. From the results, it is reasonable that local CLM feature than global salM$^3$LBP feature does much more contributions to the classification because of local details in this aerial scene, hence a smaller weight of 0.1 in global salM$^3$LBP feature and 0.9 in CLM is achieved, which results in the best accuracy according to the cross-validation searching. Visualization of confusion matrix with different classes on 30-class test set is shown in Fig. 9. Based on a visual inspection, almost all of the class-specific accuracies for our proposed method are improved, and hold for interpreting the consistency of salM$^3$LBP–CLM. All the results show that our method is very powerful for aerial scene classification.

Although different features give different performance on different scene datasets, on one hand, which can be explained by the characteristics of the dataset; on the other hand, we can observe the consistent improvements by the proposed salM$^3$LBP–CLM method on three considered scene datasets that indicates the proposed salM$^3$LBP–CLM is very effective for high-resolution scene classification. Furthermore, the results on three challenging image scene datasets demonstrate our salM$^3$LBP–CLM also can handle images with complex surroundings, such as heavy background clutters and occlusion.

### E. Computational Complexity Analysis

The computation time of the proposed approach is listed in Table V, which is simple to implement (using a PC with Intel

### TABLE V
COMPUTATION TIME (IN SECONDS) OF THE PROPOSED METHOD FOR ALL THE EXPERIMENTAL DATA

| Methods | 21-class land-use | 19-class satellite | 30-class aerial |
|---|---|---|---|
| Proposed salM$^3$LBP–CLM | 1650.5 | 7740.6 | 60062.8 |
| CLM (eSIFT) | 1302.1 | 5423.5 | 52016.4 |
| MS-CLBP+FV | 1871.4 | 8619.3 | 70101.5 |
| BoVW | 1172.2 | 5020.4 | 40033.6 |

Core i7, 32GB RAM, Windows 10, and MATLAB R2016a). For example, for the global feature extraction, the salM$^3$LBP needs about 0.35 s for images not large than $300 \times 300$ pixels and 3 s for an image of size $600 \times 600$ to build the feature histogram, whereas the local feature extracting time needs about 0.15 s for images not large than $300 \times 300$ pixels and almost 2.6 s per image with $600 \times 600$ pixels. It is noticed that global feature extraction has more complexity than local feature extraction because of a single patch size and single Gaussian modeling used in CLM, while the global features computing in multiple sampling modalities (multiscale, multiresolution, and multistructure) and stacking. According to our experiments, the proposed salM$^3$LBP–CLM needs to compute both global and local features, resulting in higher time complexity. Furthermore, the proposed salM$^3$LBP–CLM is feasible to achieve faster computation by graphical processing unit; in doing so, computing complexity will have little side-effect. Therefore, the proposed approach is effective and efficient for high-resolution scene classification.

### F. Discussion

The design of proper scene classification framework for effectively understanding the semantic content of image scenes is the first important issue we are facing due to the drastically increasing number of satellite and aerial images, the high intraclass diversity, and low interclass variations in complex scene. The proposed salM$^3$LBP–CLM exploits global salM$^3$LBP and local CLM, then fused as a final representation by cross-validation searching. An important observation is that CLM is a little sensitive to local descriptors; however, CLM is more effective than FV and BoVW. Meanwhile, the global salM$^3$LBP is proposed as complementary to CLM, which leads to substantial improvements in performance while does not increase additional computation time during the testing.

Overall, by comparing the classification performances of Sections IV-B~IV-D, it is clear that the proposed salM$^3$LBP–CLM approach is comparable or superior to the competitors in terms of classification accuracy, this is expected. The improvements mainly come from the discriminative image representation where the global salM$^3$LBP and local CLM are fused as scene descriptor, which confirms our former statement. It is interesting to note that for some difficult classes, such as sparse residential, medium residential, dense residential, and tennis court in 21-class land-use scene; industrial, residential, and parking in 19-class satellite scene; sparse residential, medium residential, and dense residential in 30-class aerial scene, no longer belong to the difficult ones. The proposed salM$^3$LBP–CLM

method is compared with the state-of-the-art such as CNN-based methods and exhibits very good generalization performance with an OA of above 90% or of 100% increase, which well validates our observation that salM$^3$LBP–CLM can improve the performance of the learnt model for difficult classes in the high-resolution image scenes. Nevertheless, as mentioned above, the proposed salM$^3$LBP–CLM is being challenged and even outperformed by recent deep CNNs. The superiority and challenge in deep learning is twofold: first, deep CNN features are learned in deep learning architectures and more powerful representations of images with multiple levels of abstraction, which leads to dramatic performance improvements for scene classification; second, it is difficult to train a deep CNN with limited training data, which typically contains millions of parameters for classification task. In total, our salM$^3$LBP–CLM achieves a good trade-off between classification accuracy and computational efficiency.

## V. Conclusion

In this paper, we propose an effective global/local feature extraction and fused representation method for high-resolution scene classification. The proposed method combines global feature based on salM$^3$LBP and local feature based on recent CLM to generate a fused representation (salM$^3$LBP–CLM) via cross-validation strategy. The proposed salM$^3$LBP–CLM can well explore the complementary attributes of structure and texture information in image scenes globally and locally, leading to substantially enhanced feature discrimination. The experimental results on three challenging scene datasets including the public largest 30-class image scene demonstrated that the proposed approach has achieved better or comparable performance as compared to the state-of-the-art methods. In future work, we will investigate hierarchical CNN features for high-resolution scene classification.

## Acknowledgment

## References

[1] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. ACM Int. Conf. Adv. Geogr. Inf. Syst.*, Nov. 2010, pp. 270–279.

[2] Q. Hu, W. B. Wu, T. Xia, Q. Y. Yu, P. Yang, Z. G. Li, and Q. Song, "Exploring the use of google earth imagery and object-based methods in land use/cover mapping," *Remote Sens.*, vol. 5, no. 11, pp. 6026–6042, 2013.

[3] A. M. Cheriyadat, "Unsupervised feature learning for aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 1, pp. 439–451, Jan. 2014.

[4] L. P. Zhang, L. F. Zhang, and B. Du, "Deep learning for remote sensing data: A technical tutorial on the state of the art," *IEEE Geosci. Remote Sens. Mag.*, vol. 4, no. 2, pp. 22–40, Jun. 2016.

[5] G. Camps-Valls, D. Tuia, L. Bruzzone, and J. A. Benediktsson, "Advances in hyperspectral image classification: Earth monitoring with statistical learning methods," *IEEE Signal Process. Mag.*, vol. 31, no. 1, pp. 45–54, Jan. 2014.

[6] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, 2004.

[7] G. F. Sheng, W. Yang, T. Xu, and H. Sun, "High-resolution satellite scene classification using a sparse coding based multiple feature combination," *Int. J. Remote Sens.*, vol. 33, no. 8, pp. 2395–2412, Apr. 2012.

[8] V. Risojevic and Z. Babić, "Aerial image classification using structural texture similarity," in *Proc. IEEE Int. Symp. Signal Process. Inf. Technol.*, 2011, pp. 190–195.

[9] V. Risojević, S. Momić, and Z. Babić, "Gabor descriptors for aerial image classification," in *Proc. Int. Conf. Adapt. Natural Comput. Algorithms*, 2011, pp. 51–60.

[10] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, no. 3, pp. 145–175, 2001.

[11] X. W. Zheng, X. Sun, K. Fu, and H. Q. Wang, "Automatic annotation of satellite images via multifeature joint sparse coding with spatial relation constraint," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 4, pp. 652–656, Jul. 2013.

[12] B. Luo, S. J. Jiang, and L. P. Zhang, "Indexing of remote sensing images with different resolutions by multiple features," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 6, no. 4, pp. 1899–1912, Aug. 2013.

[13] T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 7, pp. 971–987, Jul. 2002.

[14] Z. H. Guo, L. Zhang, and D. Zhang, "A completed modeling of local binary pattern operator for texture classification," *IEEE Trans. Image Process.*, vol. 19, no. 6, pp. 1657–1663, Jun. 2010.

[15] C. Chen, B. C. Zhang, H. J. Su, W. Li, and L. Wang, "Land-use scene classification using multi-scale completed local binary patterns," *Signal, Image Video Process.*, vol. 10, no. 4, pp. 745–752, Apr. 2016.

[16] X. Y. Bian, C. Chen, Q. Du, and Y. X. Sheng, "Extended multi-structure local binary pattern for high-resolution image scene classification," in *Proc. IEEE 36th Int. Conf. Geosci. Remote Sens. Symp.*, 2016, pp. 5134–5137.

[17] J. Yang, Y. G. Jiang, A. G. Hauptmann, and C. W. Ngo, "Evaluating bag-of-visual-words representations in scene classification," in *Proc. Int. Workshop Multimedia Inf. Retrieval*, Sep. 2007, pp. 197–206.

[18] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: spatial pyramid matching for recognizing natural scene categories," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2006, pp. 2169–2178.

[19] Y. Yang and S. Newsam, "Spatial pyramid co-occurrence for image classification," in *Proc. 2011 Int. Conf. Comput. Vis.*, Nov. 2011, pp. 1465–1472.

[20] L. Zhou, Z. T. Zhou, and D. W. Hu, "Scene classification using a multi-resolution bag-of-features model," *Pattern Recog.*, vol. 46, no. 1, pp. 424–433, Jan. 2013.

[21] L. J. Zhao, P. Tang, and L. Z. Huo, "Land-use scene classification using a concentric circle-structured multiscale bag-of-visual-words model," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 7, no. 12 pp. 4620–4631, Aug. 2014.

[22] G. Cheng, J. W. Han, L. Guo, Z. B. Liu, S. H. Bu, and J. C. Ren, "Effective and efficient midlevel visual elements-oriented land-use classification using VHR remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 8, pp. 4238–4249, Aug. 2015.

[23] R. Negrel, D. Picard, and P. H. Gosselin, "Evaluation of second-order visual features for land-use classification," in *Proc. Int. Workshop Content-Based Multimedia Indexing*, Jun. 2014, pp. 1–5.

[24] F. Perronnin, J. Sánchez, and T. Mensink, "Improving the fisher kernel for large-scale image classification," in *Proc. Eur. Conf. Comput. Vis.*, Sep. 2010, pp. 143–156.

[25] F. Zhang, B. Du, and L. P. Zhang, "Saliency-guided unsupervised feature learning for scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 4, pp. 2175–2184, Apr. 2015.

[26] O. A. B. Penatti, K. Nogueira, and J. A. dos Santos, "Do deep features generalize from everyday objects to remote sensing and aerial scenes domains?" in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2015, pp. 44–51.

[27] F. Zhang, B. Du, and L. P. Zhang, "Scene classification via a gradient boosting random convolutional network framework," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 3, pp. 1793–1802, Mar. 2016.

[28] K. Nogueira, O. A. B. Penatti, and J. A. dos Santos, "Towards better exploiting convolutional neural networks for remote sensing scene classification," *Pattern Recog.*, vol. 71, pp. 539–556, Jan. 2017.

[29] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," in *Proc. Int. Conf. Learn. Represent.*, arXiv:1312.6229, Apr. 2014.

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

BIAN *et al.*: FUSING LOCAL AND GLOBAL FEATURES FOR HIGH-RESOLUTION SCENE CLASSIFICATION

13

[30] Y. Q. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proc. ACM Int. Conf. Multimedia*, 2014, pp. 675–678.

[31] M. Castelluccio, G. Poggi, C. Sansone, and L. Verdoliva, "Land use classification in remote sensing images by convolutional networks," arXiv preprint arXiv: 1508.00092, Aug. 2015.

[32] C. Szegedy, W. Liu, Y. Q. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Boston, MA, USA, Jun. 2015, pp. 1–9.

[33] F. Hu, G. S. Xia, J. W. Hu, and L. P. 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.

[34] H. Jégou, F. Perronnin, M. Douze, J. Sanchez, P. Perez, and C. Schmid, "Aggregating local image descriptors into compact codes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 9, pp. 1704–1716, Dec. 2012.

[35] X. W. Yao, J. W. Han, G. Cheng, X. M. Qian, and L. Guo, "Semantic annotation of high-resolution satellite images via weakly supervised learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 6, pp. 3660–3671, Jun. 2016.

[36] X. Q. Lu, X. L. Li, and L. C. Mou, "Semi-supervised multi-task learning for scene recognition," *IEEE Trans. Cybern.*, vol. 45, no. 9, pp. 1967–1976, Sep. 2015.

[37] J. W. Han, D. W. Zhang, G. Cheng, L. Guo, and J. C. Ren, "Object detection in optical remote sensing images based on weakly supervised learning and high-level feature learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 6, pp. 3325–3337, Jun. 2015.

[38] G. Cheng, P. C. Zhou, and J. W. Han, "Learning rotation-invariant convolutional neural networks for object detection in VHR optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 12, pp. 7405–7415, Dec. 2016.

[39] G. Cheng and J. W. Han, "A survey on object detection in optical remote sensing images," *ISPRS J. Photogramm. Remote Sens.*, vol. 117, pp. 11–28, Jul. 2016.

[40] J. W. Han, P. C. Zhou, D. W. Zhang, G. Cheng, L. Guo, Z. B. Liu, S. H. Bu, and J. Wu, "Efficient, simultaneous detection of multi-class geospatial targets based on visual saliency modeling and discriminative learning of sparse coding," *ISPRS J. Photogramm. Remote Sens.*, vol. 89, pp. 37–48, Mar. 2014.

[41] J. Carreira, R. Caseiro, J. Batista, and C. Sminchisescu, "Freeform region description with second-order pooling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 6, pp. 1177–1189, Jun. 2015.

[42] Q. L. Wang, P. H. Li, L. Zhang, and W. M. Zuo, "Towards effective codebookless model for image classification," *Pattern Recog.*, vol. 59, pp. 63–71, Nov. 2016.

[43] C. Stein, "Lectures on the theory of estimation of many parameters," *J. Math. Sci.*, vol. 34, no. 1, pp. 1373–1403, Jul. 1986.

[44] J. Y. Zou, W. Li, C. Chen, and Q. Du, "Scene classification using local and global features with collaborative representation fusion," *Inf. Sci.*, vol. 348, pp. 209–226, Jun. 2016.

[45] G. B. Huang, H. M. Zhou, X. J. Ding, and R. Zhang, "Extreme learning machine for regression and multiclass classification," *IEEE Trans. Syst. Man Cybern.*, vol. 42, no. 2, pp. 513–529, Apr. 2012.

[46] Z. Z. Zheng, T. X. Zhang, and L. X. Yan, "Saliency model for object detection: searching for novel items in the scene," *Opt. Lett.*, vol. 37, no. 9, pp. 1580–1582, May 2012.

[47] L. Liu, L. J. Zhao, Y. L. Long, G. Y. Kuang, and P. Fieguth, "Extended local binary patterns for texture classification," *Image Vis. Comput.*, vol. 30, no. 2, pp. 86–99, Feb. 2012.

[48] L. H. Huang, C. Chen, W. Li, and Q. Du, "Remote sensing image scene classification using multi-scale completed local binary patterns and fisher vectors," *Remote Sens.*, vol. 8, no. 6, pp. 1–17, Jun. 2016.

[49] G. S. Xia, J. W. Hu, F. Hu, B. G. Shi, X. Bai, Y. F. Zhong, and L. P. Zhang, "AID: A benchmark dataset for performance evaluation of aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, to be published, 2017.

[50] W. Shao, W. Yang, G. S. Xia, and G. Liu, "A hierarchical scheme of multiple feature fusion for high-resolution satellite scene categorization," in *Proc. Int. Conf. Comput. Vis. Syst.*, Jul. 2013, pp. 324–333.

[51] S. Z. Chen and Y. L. Tian, "Pyramid of spatial relations for scene-level land use classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 4, 1947–1957, Apr. 2015.

**Xiaoyong Bian** received the B.S. degree in computer science from the Kunming University of Science and Technology, Kunming, China, in 1999, the M.S. degree in computer science from the Wuhan University of Science and Technology, Wuhan, China, in 2005, and the Ph.D. degree in pattern recognition and intelligent system from the Huazhong University of Science and Technology, Wuhan, China, in 2013.

He is currently an Associate Professor of computer science and technology with the Wuhan University of Science and Technology. During 2015, he was a Visiting Fellow in the Department of Electrical and Computer Engineering, Mississippi State University, Starkville, MS, USA. His research interests include machine learning, satellite/aerial image scene classification, and recognition.

**Chen Chen** received the B.E. degree in automation from the Beijing Forestry University, Beijing, China, in 2009, and the M.S. degree in electrical engineering from the Mississippi State University, Starkville, MS, USA, in 2012, and the Ph.D. degree in the Department of Electrical Engineering, University of Texas at Dallas, Richardson, TX, USA, in 2016.

He is currently a Post-Doc in the Center for Research in Computer Vision, University of Central Florida, Orlando, FL, USA. He has published more than 50 papers in refereed journals and conferences in these areas. His research interests include compressed sensing, signal and image processing, pattern recognition, and computer vision.

**Long Tian** received the M.S. degree in blur image detection from the Department of Electrical Engineering, City College of New York, New York, NY, USA, 2014. He is currently working toward the Ph.D. degree in the Department of Electrical and Computer Engineering, Mississippi State University, Starkville, MS, USA.

His research interests include image detection and recognition, and hyper-spectral image classification.

**Qian Du** (S'98–M'00–SM'05) received the Ph.D. degree in electrical engineering from the University of Maryland at Baltimore County, Baltimore, MD, USA, in 2000.

Currently, she is the Bobby Shackouls Professor in the Department of Electrical and Computer Engineering, Mississippi State University, Starkville, MS, USA. Her research interests include hyperspectral remote sensing image analysis and applications, pattern classification, data compression, and neural networks.

Dr. Du is a Fellow of the SPIE—International Society for Optics and Photonics. She received the 2010 Best Reviewer Award from the IEEE Geoscience and Remote Sensing Society. She was a Co-Chair for the Data Fusion Technical Committee of the IEEE Geoscience and Remote Sensing Society from 2009 to 2013, and the Chair for Remote Sensing and Mapping Technical Committee of the International Association for Pattern Recognition from 2010 to 2014. She served as an Associate Editor for the IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING, the *Journal of Applied Remote Sensing*, and the IEEE SIGNAL PROCESSING LETTERS. Since 2016, she has been the Editor-in-Chief for the IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING. She is the General Chair for the 4th IEEE GRSS Workshop on Hyperspectral Image and Signal Processing: Evolution in Remote Sensing in Shanghai, China, in 2012.