

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/268690576>

# Multi-class geospatial object detection and geographic image classification based on collection of part detectors

Article in ISPRS Journal of Photogrammetry and Remote Sensing · November 2014

DOI: 10.1016/j.isprsjprs.2014.10.002

---

CITATIONS

80

READS

445

4 authors:



Gong Cheng

Northwestern Polytechnical University

41 PUBLICATIONS 939 CITATIONS

[SEE PROFILE](#)



Junwei Han

Northwestern Polytechnical University

193 PUBLICATIONS 2,639 CITATIONS

[SEE PROFILE](#)



Peicheng Zhou

University of Technology Sydney

14 PUBLICATIONS 382 CITATIONS

[SEE PROFILE](#)



Kaiming Li

Emory University

275 PUBLICATIONS 3,509 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Remote Sensing Image Scene Classification: Benchmark and State of the Art [View project](#)



Object detection in optical remote sensing images [View project](#)





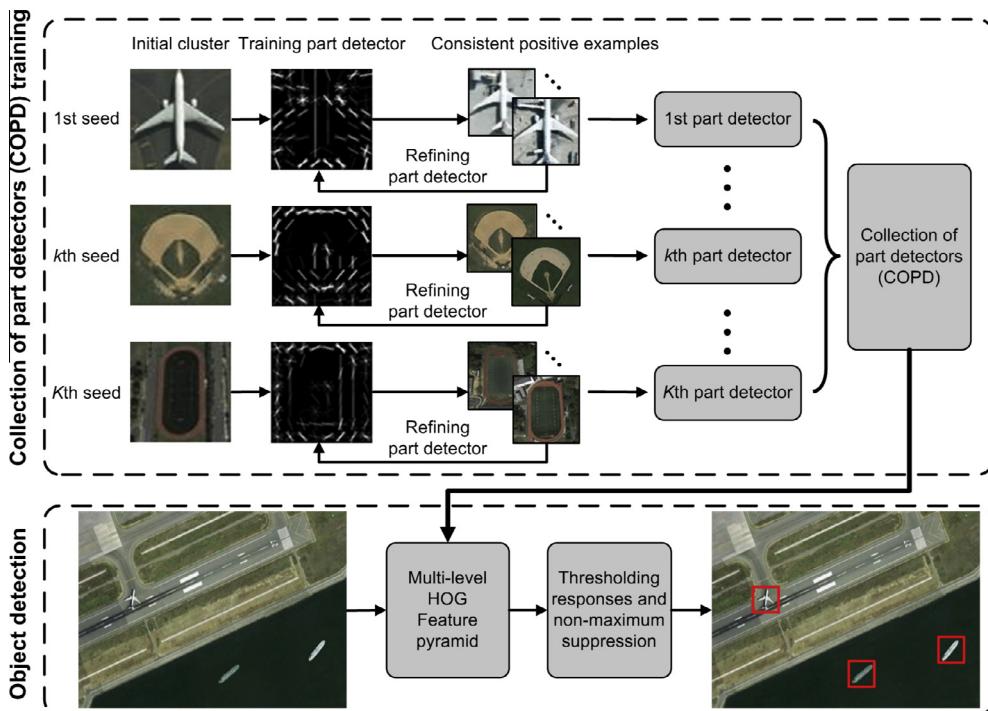


Fig. 2. Overview of the developed COPD-based multi-class geospatial object detection framework.

To sum up, the primary contribution of this paper is fourfold. First, by extending the notion of “part detector” to high-resolution remote sensing images analysis, we introduce a practical and rotation-invariant framework for multi-class geospatial object detection and geographic image classification based on collection of part detectors, where each part detector is used for the detection of objects or recurring spatial patterns within a certain range of orientation. Second, when training part detectors for multi-class object detection, we improve the traditionally used exemplar-SVM detector (Malisiewicz et al., 2011) training process by alternatively refining part detectors and incorporating consistent positive examples for each exemplar. The quantitative comparison results on first-of-its-kind 10-class objects data set, as shown in Fig. 7 and Table 3, demonstrate huge performance gain of our method compared with state-of-the-art approaches. Third, by taking advantage of the technology of mid-level visual elements discovery (Juneja et al., 2013; Li et al., 2013; Singh et al., 2012; Sun and Ponce, 2013), we achieve an effective image representation method for VHR geographic image classification by using discriminative visual parts as attributes, which provides a more informative description of an image. As shown in Fig. 12 and Table 5, superior and encouraging results are obtained on a publicly available 21-class LULC benchmark for image classification. To the best of our knowledge, this result is the best on this data set, which adequately shows the superiority and effectiveness of the developed framework. Fourth, a high-spatial-resolution remote sensing images data set containing 10-class objects (airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, and vehicle) is constructed and will be made publicly available to other researchers.<sup>1</sup> We anticipate this data set will help other researchers to conduct further study or compare different algorithms.

The rest of the paper is organized as follows. Section 2 describes a COPD-based multi-class geospatial object detection framework and reports comparative experimental results on a high resolution

remote sensing image data set. Section 3 details a COPD-based geographic image classification framework and gives comparative experimental results on a publicly available 21-class LULC benchmark. Finally, conclusions are drawn in Section 4.

## 2. COPD-based multi-class geospatial object detection

### 2.1. Framework overview

Fig. 2 gives an overview of the developed COPD-based multi-class geospatial object detection framework. It is mainly composed of two stages: COPD training and object detection. In the COPD training phase, we first pick a set of representative seeds to serve as initial clusters, where each seed corresponds to a particular viewpoint of an object class and each cluster corresponds to a part detector needed to be trained. Then, we train a set of part detectors using an iterative procedure (Bourdev and Malik, 2009; Cheng et al., 2013b; Felzenszwalb et al., 2010; Singh et al., 2012) that alternates between refining detectors and incorporating consistent positive examples for each seed from training images. Given  $K$  seeds, we can finally obtain a COPD that is composed of  $K$  seed-based part detectors. Since each part detector corresponds to a particular viewpoint of an object class, the collection of them could provide an effective solution for rotation-invariant and simultaneous detection of multi-class geospatial objects. In the object detection stage, given a new test image, we first run all detectors simultaneously on the input image, in Histograms of Oriented Gradients (HOG) (Dalal and Triggs, 2005) feature pyramid space, to obtain the response and potential object class for each sliding-window. Then, multi-class object detection is implemented by thresholding the responses and eliminating repeated detections via non-maximum suppression (Bourdev and Malik, 2009; Felzenszwalb et al., 2010).

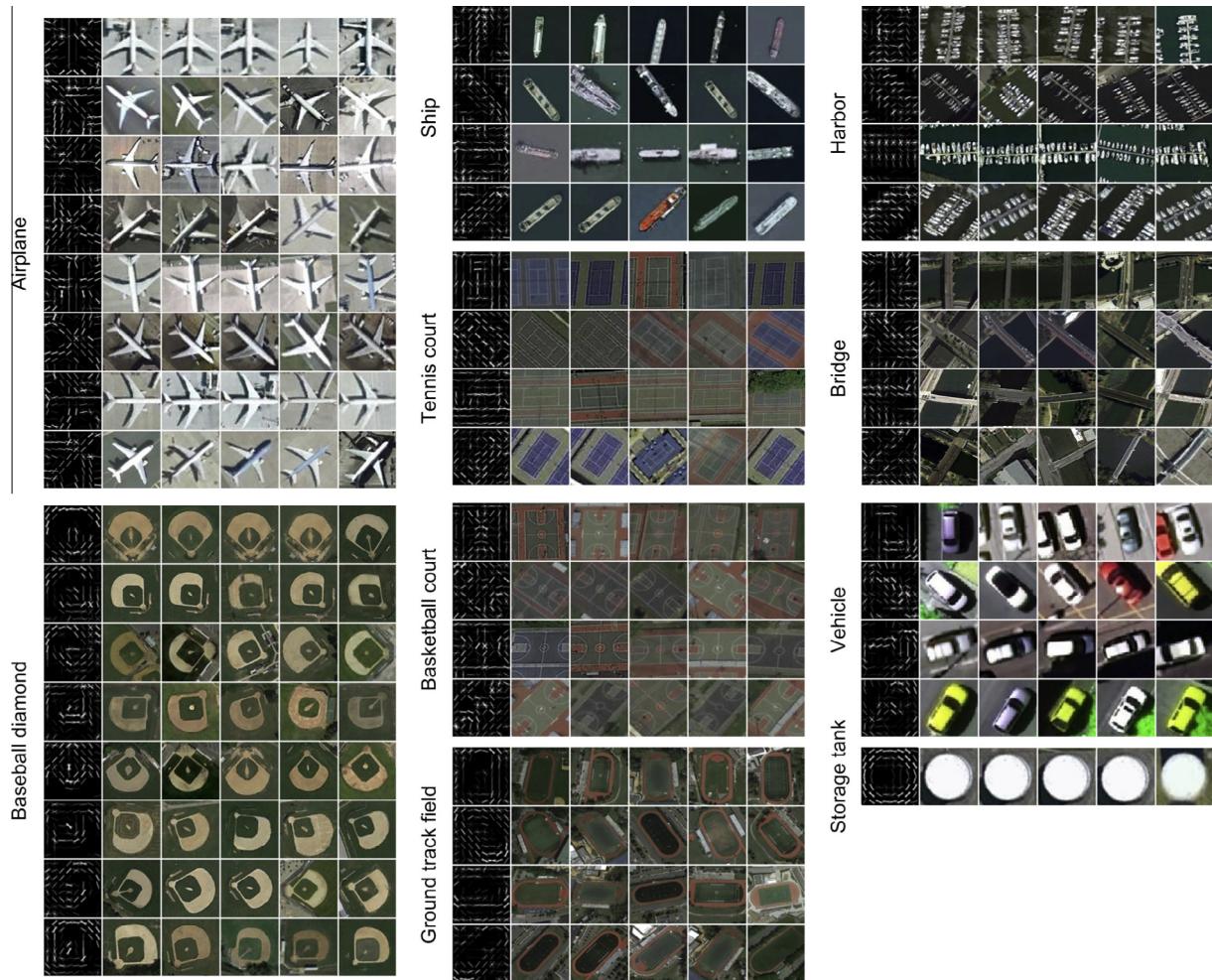
### 2.2. Framework details

#### 2.2.1. COPD training

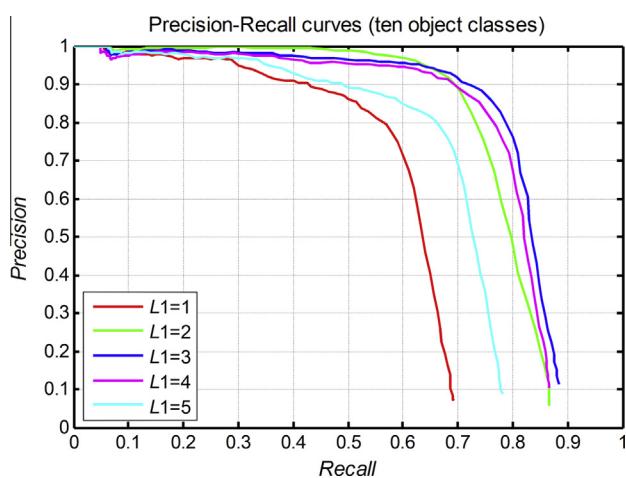
When training a COPD for multi-class object detection, the input is composed of a “positive image dataset”  $P$  in which each

<sup>1</sup> <http://pan.baidu.com/s/1c0w8h3q>.





**Fig. 4.** The visualization of weight vectors of 45 detectors for airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, and vehicle classes, respectively. Their top-5 high-scoring positives from the training data set are also shown subsequently. We have resized these positives to  $60 \times 60$  pixels for visualization.



**Fig. 5.** Precision–Recall curves obtained by varying the values of  $L_1$ .

**Table 2**  
Performance comparisons of different  $L_1$  in terms of AP.

$L_1$	1	2	3	4	5
AP	0.5929	0.7798	0.8044	0.7838	0.6798

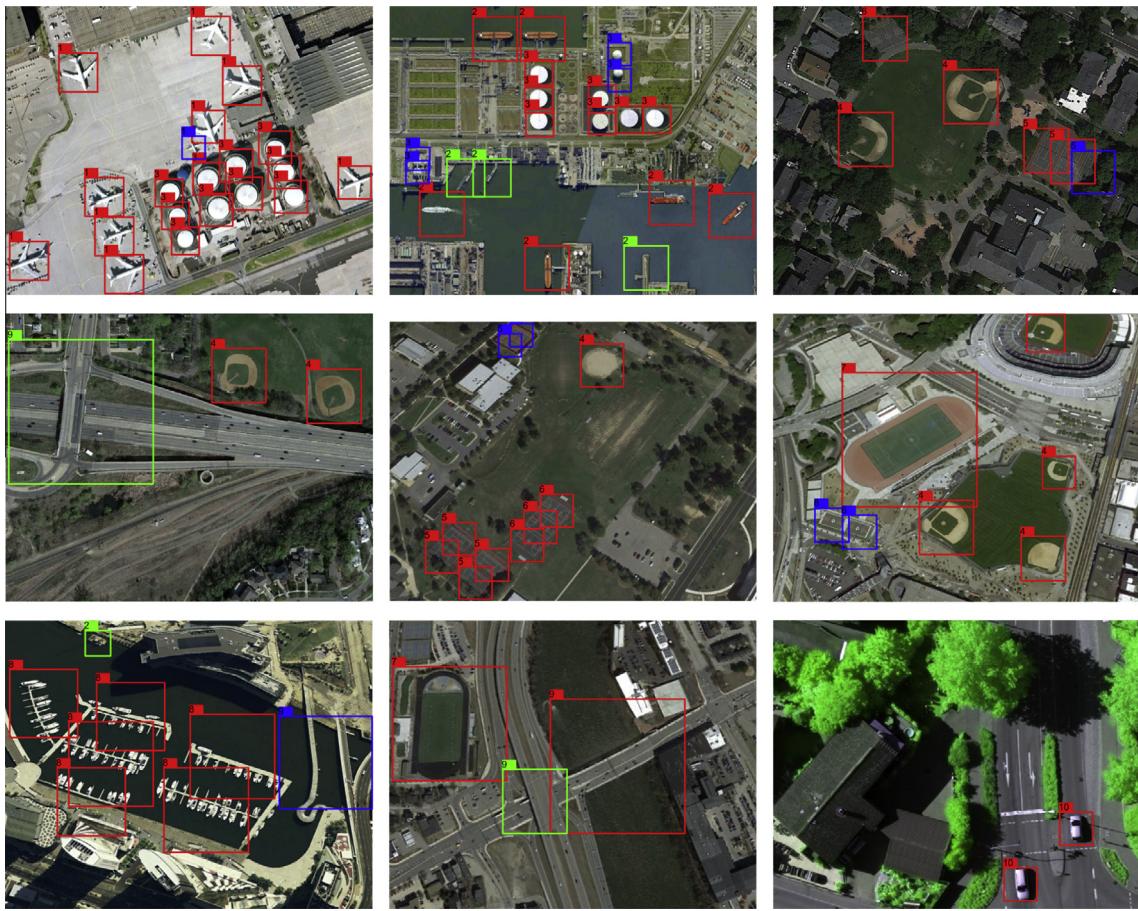
select the highest scoring ones while removing those that are at least 50% covered by a previously selected bounding box.

### 2.3. Experiments

#### 2.3.1. Data set description

In theory, the developed multi-class object detection framework can detect a large number of classes of geospatial objects. However, in our experiments, we used the task of detection of ten different types of objects to evaluate the performance of the developed framework. These ten classes of objects are airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, and vehicle.

We collected 715 high-spatial-resolution color images from Google Earth and 85 very-high-spatial-resolution pansharpened color infrared (CIR) images from Vaihingen data set (Cramer, 2010) used for our evaluations, where the spatial resolution of Google Earth images ranges from 0.5 m to 2 m and the spatial resolution of CIR images is 0.08 m. The Vaihingen data was provided by the German Society for Photogrammetry, Remote Sensing and Geoinformation (DGPF) (Cramer, 2010): <http://www.ifp.uni-stuttgart.de/dgpf/DKEPAllg.html>. We divided these images into four independent datasets: a “negative image set” containing 150 images, a “positive image set” containing 150 images, an “optimizing set” containing 150 images, and a testing set containing 350



**Fig. 6.** A number of multi-class object detection results by using the developed framework.

images. All images from the first set do not contain any targets of the given object classes and each image from the last three sets contains at least one target to be detected. The “negative image set” and “positive image set” were used for the COPD training, the “optimizing set” was used for parameter optimization and the “testing set” was used for testing the performance of the developed framework. We labeled ground truths from the “optimizing set” and the “testing set”, respectively. The detailed object sizes, object numbers from optimizing set, and object numbers from testing set of ten different object classes are listed in Table 1.

### 2.3.2. Seeds generation

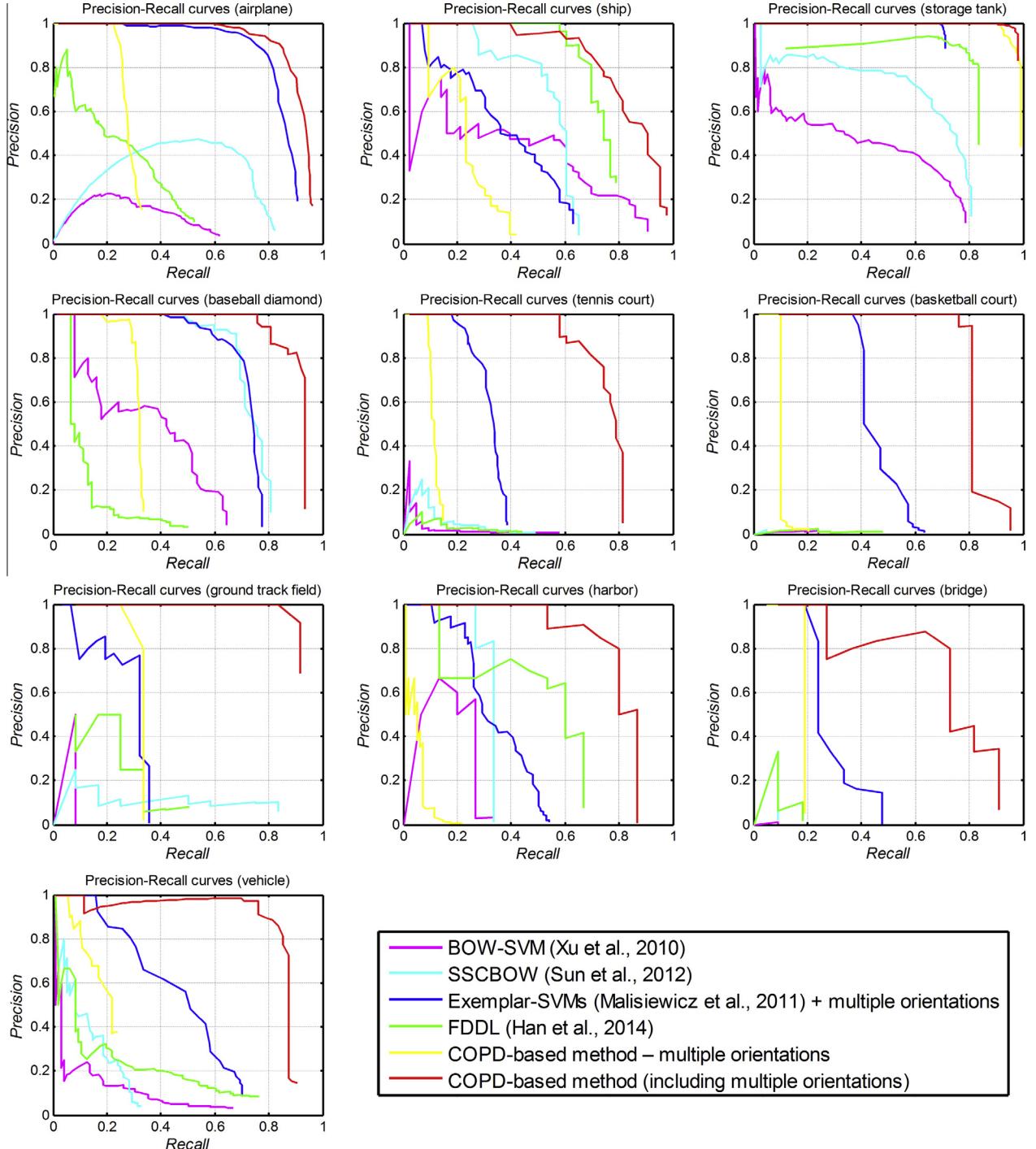
Here, it should be pointed out that the seeds serving as initial clusters should be representative and have different orientations, which can be obtained by manually labeling a representative sample for each object class from the “positive image set”, aligning them, and then rotating them with a certain angle. Specifically, for our 10-class object detection task, given ten manually labeled samples (each sample for each object class), we first align each of them to an unified orientation (e.g. approximately vertical in our implementation) and then rotate airplane and baseball diamond samples in the step of  $45^\circ$  from  $0^\circ$  to  $360^\circ$ , rotate samples of ship, tennis court, basketball court, ground track field, harbor, bridge, and vehicle in the step of  $45^\circ$  from  $0^\circ$  to  $180^\circ$  because their shapes are bilaterally symmetric, and perform no rotation for storage tank sample because its shape is circular. In this way, we can obtain eight seeds for each manually labeled sample of airplane and baseball diamond, four seeds for each manually labeled sample of ship, tennis court, basketball court, ground track field, harbor, bridge, and vehicle, and one seed for each manually labeled sample of

storage tank. Fig. 3 illustrates the total 45 seeds used in our work for 10-class object detection.

Using the COPD training procedure as described in Subsection 2.2.1 and the 45 seeds as shown in Fig. 3, we trained a COPD consisting of 45 detectors for airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, and vehicle classes on our training data set. Fig. 4 shows the visualization of weight vectors  $w_k$  of 45 detectors, where brighter “pixel” represents bigger weight and vice versa. Their corresponding top-5 high-scoring positives from the training data set are also shown subsequently.

### 2.3.3. Evaluation criterions

We consider a detection to be correct if its bounding box overlaps more than 50% with the ground truth bounding box, otherwise the detection is considered as a false positive. In addition, if several bounding boxes overlap with a same single ground truth bounding box, only one is considered as true positive and the others are considered as false positives. We adopted the standard Precision–Recall curve (PRC) (Buckland and Gey, 1994) and AP (Everingham et al., 2010) to quantitatively evaluate the performance of an object detection system. The *Precision* measures the fraction of detections that are true positives and the *Recall* measures the fraction of positives that are correctly identified. AP computes the average value of *Precision* over the interval from  $Recall = 0$  to  $Recall = 1$ , i.e. the area under the PRC, so the higher the AP value is, the better the performance and vice versa. Let  $TP$ ,  $FP$ , and  $NP$  denote the number of true positives, the number of false positives, and the number of total positives. The *Precision* and *Recall* can be formulated as:



**Fig. 7.** Precision–Recall curves of the developed framework and some state-of-the-art approaches for airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, and vehicle classes respectively.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (3)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{NP}} \quad (4)$$

#### 2.3.4. Experimental setting

In the implementation of multi-class object detection, to address the problem that the sizes of targets may be different in images, each image is represented by an 15-level HOG (Dalal and Triggs, 2005) feature pyramid and each octave contains five levels (i.e. for  $l$  th level, the sub-sampling factor is  $2^{(l-1)/5}$ ). We follow the

construction in Dalal and Triggs (2005) to extract the HOG feature for each pyramid level. Specifically, we partition the image at each pyramid level into non-overlapping cells of  $6 \times 6$  pixels and use nine orientation bins to accumulate a one-dimensional histogram of gradient orientations over pixels in each cell. Then, each  $2 \times 2$  neighbourhood of cells is grouped into one block (with a stride of one cell) and a robust normalization process based on 2-norm is run on each block to provide greater invariance to local illumination and spatial deformation, which finally forms a 36-dimensional HOG feature vector. Rather than using the 36-dimensional vector directly, in this work we project it onto a lower 31-dimensional



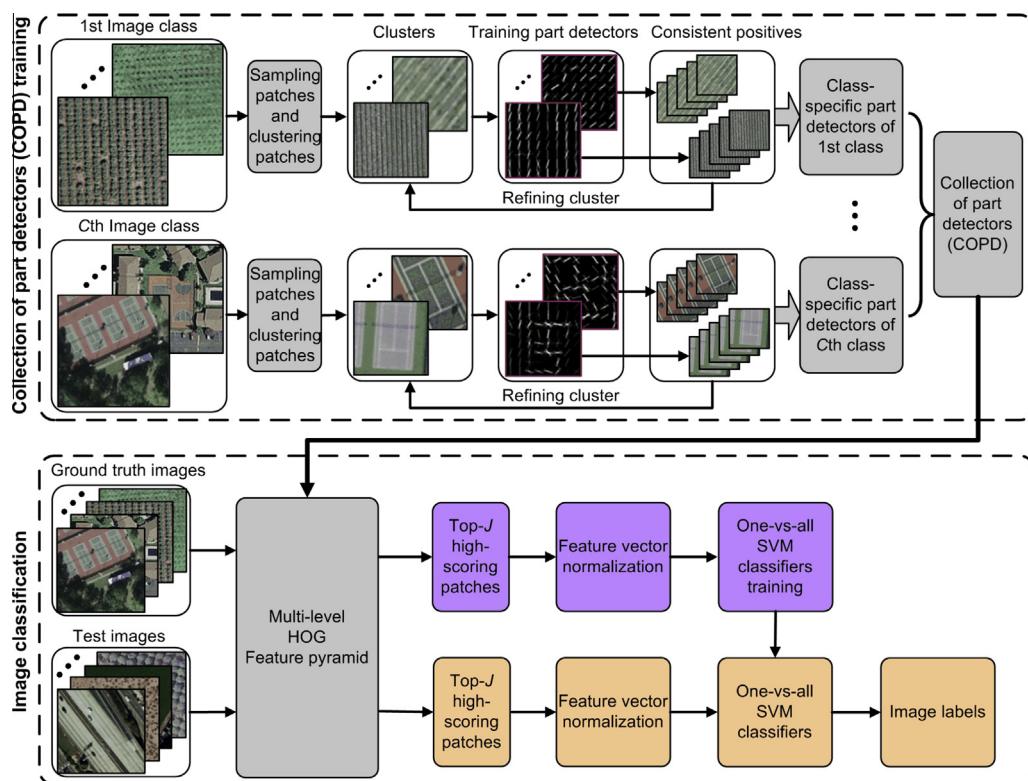
diamond targets; (2) Similar to BOW-SVM method, SSCBOW method (Sun et al., 2012) also represents each image patch as a histogram of visual words, but in which sparse coding is introduced to replace  $K$ -means algorithm for visual words encoding. This new spatial encoding strategy not only represents the relative position of the local features but also has the ability to encode the geometric information of an object, so SSCBOW method obtained better performance compared with BOW-SVM method. However, the detection results of these two methods depend largely on the extracted keypoints or local features such as SIFT descriptors. For those objects (e.g. tennis courts and basketball courts characterized by side lines, center lines, and goal lines) from which enough and discriminative keypoints are difficult to extract, the detection performances of these two methods are severely limited; (3) The method of ‘Exemplar-SVMs’ (Malisiewicz et al., 2011) + multiple orientations’ is based on training an individual linear SVM classifier for every selected exemplar in HOG feature space. Since each of these Exemplar-SVMs is defined by a single positive instance, each detector is quite specific to its exemplar and has poor generalization because of the appearances variation and deformation of objects; (4) In FDDL method, sparse representation based classification strategy is adopted to perform multi-class object detection, in which each image patch is described by a few representative atoms of a learned dictionary in a low-dimensional manifold. Unfortunately, as image patches need to be down-sampled to adapt the size of atoms, some critical and discriminative features of them (e.g. the straight lines and arc in tennis courts and basketball courts) are reduced and even removed. This fatal operation has significantly degenerated the detection accuracy; (5) Our developed framework performs object detection using a collection of part detectors derived from a set of representative seeds, where each detector corresponds to a particular viewpoint of an object class and is trained using an iterative procedure that alternatively refines part detectors and incorporates consistent positive examples for each seed from the training images. On the one hand,

incorporation of consistent positive examples could guarantee the learned detectors have good generalization and therefore can effectively handle object deformations and appearance variations compared to traditionally used Exemplar-SVMs (Malisiewicz et al., 2011). On the other hand, since each part detector corresponds to a particular viewpoint of an object class, the collection of them could provide an effective solution for rotation-invariant and simultaneous detection of multi-class geospatial objects. Consequently, our method can obtain promising results compared to the aforementioned four state-of-the-art approaches.

### 3. COPD-based geographic image classification

#### 3.1. Framework overview

The flowchart of our COPD-based geographic image classification is illustrated in Fig. 8. It is mainly composed of two stages: **COPD training** and **image classification**. In the COPD training stage, given an image database, we first train class-specific part detectors for each image class based on class labels in a weakly supervised fashion. This can be achieved by sampling a large number of image patches from the positive training images, clustering them, and alternating between training discriminative part detectors and refining clusters (Bourdev and Malik, 2009; Felzenszwalb et al., 2010; Singh et al., 2012). Then, all part detectors of each image class are combined to generate a complete COPD. In the image classification stage, we first use the trained COPD to detect mid-level visual elements (i.e. discriminative image patches) from each image and represent the image as a feature vector of the responses of the top- $J$  high-scoring image patches, which can provide more informative description of the image. Then, we train a linear one-vs-all SVM classifier for each image class by treating the images of the chosen class as positive instances and the rest images as negative instances. Finally, a test image is assigned to the label of the classifier with the highest response.



**Fig. 8.** Flowchart of the developed COPD-based geographic image classification framework.

### 3.2. Framework details

#### 3.2.1. COPD training

Let  $\Omega = \{\omega_1, \omega_2, \dots, \omega_C\}$  denote the set of  $C$  image classes of an image database. We first learn class-specific part detectors  $\Gamma^c = \{\Gamma_k^c\}_{k=1}^{K_c}$  for each image class  $\omega_c (c = 1, \dots, C)$ , where  $K_c$  is the total number of part detectors and  $\Gamma^c$  should be representative for class  $\omega_c$  and be discriminative against classes ( $\Omega - \omega_c$ ). Then, all part detectors of each class are combined to generate a complete COPD  $\Gamma = \{\Gamma^c\}_{c=1}^C$ . The training of  $\Gamma^c$  for class  $\omega_c$  is performed in terms of the following steps (Felzenszwalb et al., 2010; Singh et al., 2012; Sun and Ponce, 2013):

- (1) Construct “positive image dataset”  $P_c$  and “negative image dataset”  $N_c$ , where  $P_c$  is composed of the images of class  $\omega_c$  and  $N_c$  is composed of the images of the classes ( $\Omega - \omega_c$ ).
- (2) Randomly crop a large number of image patches from all images in  $P_c$  at different image scales, discard highly overlapping patches, perform standard k-means clustering over these image patches in HOG (Dalal and Triggs, 2005) feature space, and then retain sufficiently large clusters with size of 10 or more. The cluster number is adaptively set to be one tenth of the total number of sampled image patches in our work.
- (3) Train a linear SVM classifier  $\Gamma_k^c = (w_k^c, b_k^c) (k = 1, \dots, K_c)$  for each cluster in HOG (Dalal and Triggs, 2005) feature space, using image patches within the cluster as positive examples and all hard negative examples of  $N_c$  as negative examples. Learning the parameters  $w_k^c$  and  $b_k^c$  amounts to optimizing the similar objective function as illustrated in Eq. (1).
- (4) Run  $\Gamma^c = \{\Gamma_k^c\}_{k=1}^{K_c}$  on  $P_c$  to form new clusters from the top- $m$  high-scoring patches for each part detector. In our

implementation, we set  $m = 10$  to keep each cluster having a high purity.

- (5) Repeat the steps of (3) and (4)  $L_2$  iterations until the maximum measured by image classification accuracy is reached, thus we can obtain a updated class-specific part detectors  $\Gamma^c$  for image class  $\omega_c$ .

In our implementation, the aforementioned training procedure comes to a maximum after 4 iterations. We will report the effect of the parameter  $L_2$  in subsection 3.3. Using the above procedure, we trained five COPDs for all five held-out sets on a publicly available 21-class LULC data set (data set will be described in subsection 3.3.1). The total numbers of part detectors  $K = \sum_{c=1}^C K_c$  on all five held-out sets are  $K = \{3093, 3140, 3072, 3110, 2822\}$ . Fig. 9 shows the visualization of two randomly selected part detectors for each image class from the first held-out set, and their corresponding top-5 high-scoring image patches. It is very interesting to see that the part detectors can capture more informative visual elements that seem very intuitive to us. For example, the part detectors for the “airplane” class capture the airplanes with different orientations and sizes; the ones for the “intersection” category capture the turnings and the zebra crossings. These discriminative detectors can therefore capture the essence of the scene in terms of these highly consistent and repeating patterns and hence provide a conceptually simple but surprisingly effective visual representation.

#### 3.2.2. COPD-based image representation

Image representation plays a key role in scene-level geographic image classification. In this work, we present an effective image

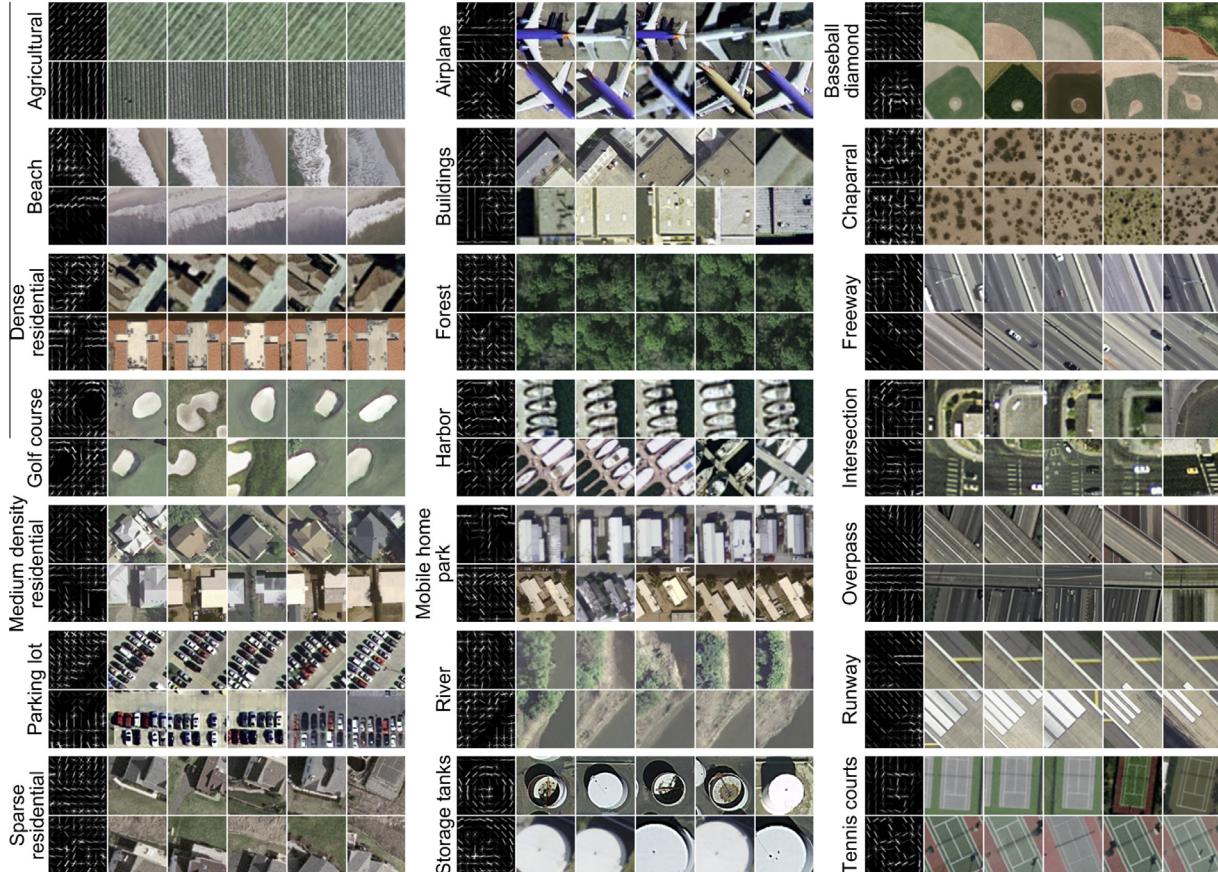
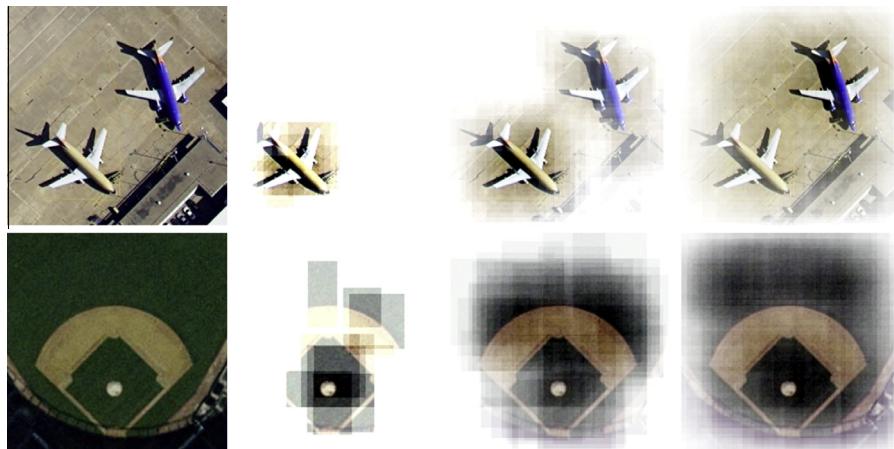


Fig. 9. The visualization of two randomly selected part detectors for each image class and their corresponding top-5 high-scoring image patches.



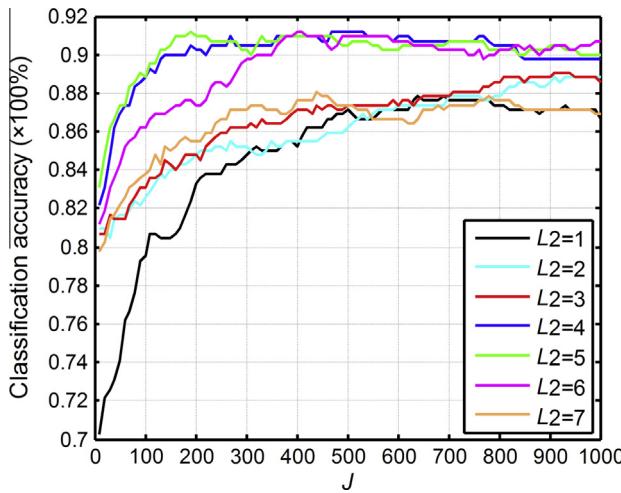
**Fig. 10.** Two original images (1st column) and their visualization images constructed by averaging their top-10 (2nd column), top-100 (3rd column), and top-500 (4th column) high-scoring patches.

representation method based on a collection of representative and discriminative part detectors that uses mid-level visual elements (i.e., discriminative image patches) as attributes for image representation. Its core is that through detecting discriminative image patches by using pre-trained part detectors, each image is represented as a feature vector of the responses of the top- $J$  high-scoring image patches. For example, Fig. 10 shows two original images (1st column) and their visualization images constructed by averaging their top-10 (2nd column), top-100 (3rd column), and top-500 (4th

column) high-scoring image patches, respectively. As can be seen from Fig. 10, a smaller value of  $J$  cannot obtain all image patches that are most related to the image class (e.g. the second column in Fig. 10). A bigger value of  $J$  can result in certain background except for the image patches that are most related to the image class (e.g. the fourth column in Fig. 10). Therefore, selecting an optimal parameter of  $J$  to capture the most discriminative essence of the scene is very important for high-level scene recognition tasks. We will report the detailed parameter optimization in Subsection 3.3.2.



**Fig. 11.** Some example images from the 21-class LULC data set.



**Fig. 12.** Classification accuracy obtained by varying  $L_2$  and  $J$ .

Specifically, given an input image  $I$ , we first run all part detectors  $\Gamma = \{\Gamma^c\}_{c=1}^C$  on its HOG (Dalal and Triggs, 2005) feature pyramid  $H(I)$  to compute the response and corresponding part detector label for each location by adopting the similar process used for object detection, as illustrated in Eq. (2). Next, the top- $J$  high-scoring image patches (measured by their responses) and their part detector labels are obtained. Finally, the responses are normalized to  $[0, 1]$  and the input image is represented as a feature vector of  $F(I)$  by accumulating all normalized responses to their corresponding part detectors. The dimension of  $F(I)$  equals to the total number of part detectors in  $\Gamma = \{\Gamma^c\}_{c=1}^C$ .

**Table 4**

Classification accuracies over all 21 classes for five different held-out sets.

Held-out set number	1	2	3	4	5	Average
Classification accuracies (%)	90.95	90.24	93.33	90.48	91.67	$91.33 \pm 1.11$

### 3.2.3. Image classification

We use a simple one-vs-all scheme to perform image classification by constructing a set of binary SVM classifiers. Each one-vs-all SVM classifier is trained individually by treating the images of the chosen class as positive instances and the rest images as negative instances. An unlabeled test image is assigned to label of the classifier with the highest response.

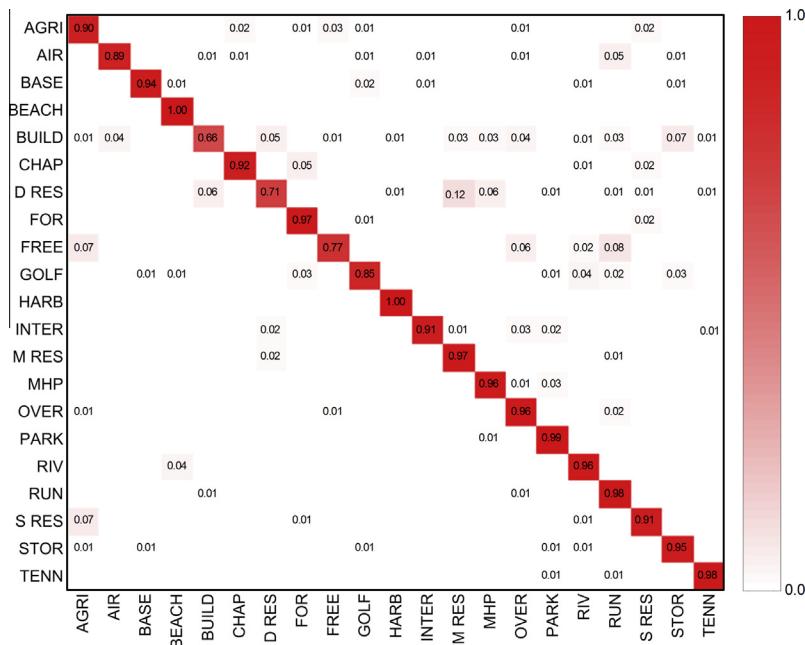
### 3.3. Experiments

#### 3.3.1. LULC data set description

We comprehensively evaluate the performance of the explored COPD-based image classification method on a publicly available data set downloaded from <http://vision.ucmerced.edu/datasets> (Yang and Newsam, 2010, 2011). The data set comprises the following 21 LULC classes: agricultural, airplane, baseball diamond, beach, buildings, chaparral, dense residential, forest, freeway, golf course, harbor, intersection, medium density residential, mobile home park, overpass, parking lot, river, runway, sparse residential, storage tanks, and tennis courts. Each class consists of 100 images measuring  $256 \times 256$  pixels, with a pixel resolution of 30 cm in the red-green-blue color space. Fig. 11 shows four samples of each class from this data set.

#### 3.3.2. Experimental setting

In the implementation of image classification, when detecting discriminative patches from images, we need construct a multi-level HOG feature pyramid for each image in a similar way as object detection. The only difference between them is that the total number of feature pyramid level  $L$  is not limited to eight and it changes as the image size changes, i.e.  $L = \lfloor 5\log_2 \min(\text{rows}, \text{cols})/60 \rfloor + 1$ , where  $\text{rows}$  and  $\text{cols}$  denote the image size in pixels in row and



**Fig. 13.** Confusion matrix by averaging classification results over all five cross-validations.



- Felzenszwalb, P.F., Girshick, R.B., McAllester, D., Ramanan, D., 2010. Object detection with discriminatively trained part-based models. *IEEE Trans. Pattern Anal. Mach. Intell.* 32, 1627–1645.
- Grabner, H., Nguyen, T.T., Gruber, B., Bischof, H., 2008. On-line boosting-based car detection from aerial images. *ISPRS J. Photogramm. Remote Sens.* 63, 382–396.
- Han, J., Zhou, P., Zhang, D., Cheng, G., Guo, L., Liu, Z., Bu, S., Wu, J., 2014. Efficient, simultaneous detection of multi-class geospatial targets based on visual saliency modeling and discriminative learning of sparse coding. *ISPRS J. Photogramm. Remote Sens.* 89, 37–48.
- Juneja, M., Vedaldi, A., Jawahar, C., Zisserman, A., 2013. Blocks that shout: distinctive parts for scene classification. In: Proceedings of the 2013 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2013), Portland, OR, pp. 923–930.
- Kim, M., Madden, M., Warner, T.A., 2009. Forest type mapping using object-specific texture measures from multispectral Ikonos imagery: segmentation quality and image classification issues. *Photogramm. Eng. Remote Sens.* 75, 819–829.
- Lazebnik, S., Schmid, C., Ponce, J., 2006. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In: Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2006). IEEE, New York, pp. 2169–2178.
- Li, F.F., Perona, P., 2005. A bayesian hierarchical model for learning natural scene categories. In: Proceedings of the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2005). IEEE, San Diego, CA, pp. 524–531.
- Li, Q., Wu, J., Tu, Z., 2013. Harvesting mid-level visual concepts from large-scale internet images. In: Proceedings of the 2013 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2013), Portland, OR, pp. 851–858.
- Lowe, D.G., 2004. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision* 60, 91–110.
- Malisiewicz, T., Gupta, A., Efros, A.A., 2011. Ensemble of exemplar-svms for object detection and beyond. In: Proceedings of the thirteenth IEEE International Conference on Computer Vision (ICCV 2011). IEEE, Barcelona, Spain, pp. 89–96.
- Martha, T.R., Kerle, N., van Westen, C.J., Jetten, V., Kumar, K.V., 2011. Segment optimization and data-driven thresholding for knowledge-based landslide detection by object-based image analysis. *IEEE Trans. Geosci. Remote Sens.* 49, 4928–4943.
- Schroder, M., Rehrauer, H., Seidel, K., Datcu, M., 2000. Interactive learning and probabilistic retrieval in remote sensing image archives. *IEEE Trans. Geosci. Remote Sens.* 38, 2288–2298.
- Shyu, C., Klaric, M., Scott, G.J., Barb, A.S., Davis, C.H., Palaniappan, K., 2007. GeoIRIS: Geospatial information retrieval and indexing system—content mining, semantics modeling, and complex queries. *IEEE Trans. Geosci. Remote Sens.* 45, 839–852.
- Singh, S., Gupta, A., Efros, A.A., 2012. Unsupervised discovery of mid-level discriminative patches. In: Proceedings of the twelfth European Conference on Computer Vision (ECCV 2012). Springer, Firenze, Italy, pp. 73–86.
- Sirmacek, B., Ünsalan, C., 2011. A probabilistic framework to detect buildings in aerial and satellite images. *IEEE Trans. Geosci. Remote Sens.* 49, 211–221.
- Sun, J., Ponce, J., 2013. Learning discriminative part detectors for image classification and cosegmentation. In: Proceedings of the fourteenth IEEE International Conference on Computer Vision (ICCV 2013), Sydney, Australia, pp. 3400–3407.
- Sun, H., Sun, X., Wang, H., Li, Y., Li, X., 2012. Automatic target detection in high-resolution remote sensing images using spatial sparse coding bag-of-words model. *IEEE Geosci. Remote Sens. Lett.* 9, 109–113.
- Ünsalan, C., Sirmacek, B., 2012. Road network detection using probabilistic and graph theoretical methods. *IEEE Trans. Geosci. Remote Sens.* 50, 4441–4453.
- Văduva, C., Gavăt, I., Datcu, M., 2013. Latent dirichlet allocation for spatial analysis of satellite images. *IEEE Trans. Geosci. Remote Sens.* 51, 2770–2786.
- Xu, S., Fang, T., Li, D., Wang, S., 2010. Object classification of aerial images with bag-of-visual words. *IEEE Geosci. Remote Sens. Lett.* 7, 366–370.
- Yang, Y., Newsam, S., 2010. Bag-of-visual-words and spatial extensions for land-use classification. In: Proceedings of the eighteenth SIGSPATIAL International Conference on Advances in Geographic Information Systems. ACM, San Jose, California, pp. 270–279.
- Yang, Y., Newsam, S., 2011. Spatial pyramid co-occurrence for image classification. In: Proceedings of the thirteenth IEEE International Conference on Computer Vision (ICCV 2011). IEEE, Barcelona, Spain, pp. 1465–1472.
- Yang, Y., Newsam, S., 2013. Geographic image retrieval using local invariant features. *IEEE Trans. Geosci. Remote Sens.* 51, 818–832.
- Zhu, C., Zhou, H., Wang, R., Guo, J., 2010. A novel hierarchical method of ship detection from spaceborne optical image based on shape and texture features. *IEEE Trans. Geosci. Remote Sens.* 48, 3446–3456.