

# Detection of compound structures using a Gaussian mixture model with spectral and spatial constraints

Çağlar Arı<sup>a</sup> and Selim Aksoy<sup>b</sup>

<sup>a</sup>Dept. of Electrical and Electronics Engineering, Bilkent University, Ankara, 06800, Turkey

<sup>b</sup>Dept. of Computer Engineering, Bilkent University, Ankara, 06800, Turkey

## ABSTRACT

High spectral and high spatial resolution images acquired from new generation satellites have enabled new applications. However, the increasing amount of detail in these images also necessitates new algorithms for automatic analysis. This paper describes a new approach to discover compound structures such as different types of residential, commercial, and industrial areas that are comprised of spatial arrangements of primitive objects such as buildings, roads, and trees. The proposed approach uses a robust Gaussian mixture model (GMM) where each Gaussian component models the spectral and shape content of a group of pixels corresponding to a primitive object. The algorithm can also incorporate spatial constraints on the layout of the primitive objects in terms of their relative positions. Given example structures of interest, a new learning algorithm fits a GMM to the image data, and this model can be used to detect other similar structures by grouping pixels that have high likelihoods of belonging to the Gaussian object models while satisfying the spatial layout constraints without any requirement for region segmentation. Experiments using WorldView-2 data show that the proposed method can detect high-level structures that cannot be modeled using traditional techniques.

**Keywords:** Object detection, Gaussian mixture model, expectation-maximization, maximum likelihood estimation, robust estimation, constrained estimation

## 1. INTRODUCTION

Recently available multispectral channels in very high spatial resolution (VHR) images acquired from new generation satellites have enabled new applications as the increased spectral resolution enhanced the capability to distinguish different physical materials. However, the increased amount of spatial detail in these images also necessitates new advanced algorithms for automatic analysis. For example, the commonly used classification algorithms that require an initial segmentation of the image into homogeneous regions cannot always cope with the increasing complexity because such homogeneous regions often correspond to very small details.

An alternative approach in the recent years has been to model the spatial arrangements of simple image regions to identify complex region groups. Gaetano et al.<sup>1</sup> performed hierarchical texture segmentation assuming that frequent neighboring regions are strongly related. They clustered the image regions to compute the frequencies of quantized region pairs with discrete labels, and used these frequencies to build a segmentation tree where some of the nodes correspond to complex structures. Zamalieva et al.<sup>2</sup> found the significant relations between neighboring regions as the modes of a probability distribution estimated using the continuous features of region co-occurrences. The resulting modes were used to construct the edges of a graph, and a graph mining algorithm was used to find subgraphs that may correspond to compound structures. Vanegas et al.<sup>3</sup> proposed a method based on fuzzy measures of relative direction between objects to detect aligned object groups. They first detected locally aligned groups of three objects, and then checked for global alignment using these local alignments. Akcay and Aksoy<sup>4</sup> described a procedure that combined statistical characteristics of primitive objects modeled using spectral, shape, and position information with structural characteristics encoded using spatial alignments of neighboring similar object groups. However, all of these approaches required an initial segmentation for the

---

Further author information: (Send correspondence to S.A.)

C.A.: E-mail: cari@ee.bilkent.edu.tr

S.A.: E-mail: saksoy@cs.bilkent.edu.tr, Telephone: +90 (312) 2903405

identification of the primitive regions. Furthermore, they were designed to detect only a particular type of arrangement such as co-occurrence or alignment.

This paper describes a new approach that combines statistical and structural characteristics of simple objects to discover compound structures in VHR images. The compound structures of interest can include different types of residential, commercial, industrial, and agricultural areas that are comprised of spatial arrangements of primitive objects such as buildings, roads, and trees corresponding to locally homogeneous details. The proposed approach uses a probabilistic representation of the image content by providing a robust extension to the commonly used Gaussian mixture models (GMM). In this model, each pixel is represented using a feature vector that encodes both spectral and spatial information consisting of the pixel's multispectral data and its coordinates, respectively. Then, each Gaussian component in the GMM models a group of pixels corresponding to a particular object where the spectral mean corresponds to the color of the object, the spectral covariance corresponds to the homogeneity of the color content, the spatial mean corresponds to the position of the object, and the spatial covariance models its shape.

Given example compound structures of interest that are comprised of multiple primitive objects, first, a GMM is fit to the pixels corresponding to the selected structures. This GMM is used as the reference model for identifying the occurrences of other similar compound structures. We describe a novel Expectation-Maximization (EM) based learning algorithm that adapts a new GMM to new image data so that the resulting Gaussian components are similar to those in the reference GMM, and the areas that do not have any similarity to these components are rejected as outliers. The algorithm also incorporates spatial constraints on the layout of the primitive objects in terms of their relative positions within the compound structure. The result is a list of compound structures detected in other parts of the same image or in other images by grouping pixels that have high likelihoods of belonging to the Gaussian object models while satisfying the spatial layout constraints. A very important feature of the proposed model is that it can perform object detection without any requirement of initial segmentation where the only assumption is that the spectral and spatial content of the primitive objects can be modeled in terms of Gaussians.

The rest of the paper is organized as follows. Section 2 defines the properties of compound structures of interest. Section 3 describes the proposed Gaussian mixture model. Section 4 presents the learning algorithm that adapts this model to the image data. Section 5 provides experimental results on an 8-band multispectral WorldView-2 image of Ankara, Turkey. Finally, Sec. 6 lists the conclusions.

## 2. COMPOUND STRUCTURES

In this paper, compound structures are defined as high-level objects with heterogeneous content that are composed of multiple primitive objects. These primitive objects typically have relatively uniform spectral content and simple shapes. We assume that the spectral and shape content of the primitive objects can be modeled using Gaussians. Consequently, we represent each pixel using a  $d$ -dimensional feature vector  $\mathbf{x} \in \mathbb{R}^d$  where  $\mathbf{x} = (\mathbf{x}^{ms}; \mathbf{x}^{xy})$  is formed by concatenating a  $d - 2$  dimensional vector  $\mathbf{x}^{ms}$  containing the multispectral values and a 2-dimensional vector  $\mathbf{x}^{xy}$  containing the pixel's coordinates. We further assume that the multispectral values and the pixel coordinates are independent, i.e.,  $p(\mathbf{x}) = p(\mathbf{x}^{ms})p(\mathbf{x}^{xy})$ . Thus, the Gaussian model for the primitive object is defined in terms of the mean  $\boldsymbol{\mu} = (\boldsymbol{\mu}^{ms}; \boldsymbol{\mu}^{xy})$  and the covariance matrix with a block diagonal structure  $\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}^{ms}, 0; 0, \boldsymbol{\Sigma}^{xy})$ . Given a group of pixels forming the primitive object, the spectral mean  $\boldsymbol{\mu}^{ms}$  corresponds to the average color of the object, the spectral covariance  $\boldsymbol{\Sigma}^{ms}$  corresponds to the homogeneity of the color content, the spatial mean  $\boldsymbol{\mu}^{xy}$  corresponds to the position of the object, and the spatial covariance  $\boldsymbol{\Sigma}^{xy}$  models its shape. We expect that these Gaussians are sufficient and effective models because of the expected relatively homogeneous spectral content and simple shapes of the primitive objects.

The primitive objects form different compound structures according to different spatial layouts. In this paper, we use a very simple layout model defined in terms of the displacement vectors between the centroids (spatial means)  $\boldsymbol{\mu}^{xy}$  of the primitive objects. Given  $K$  primitive objects, the spatial layout of the compound structure is modeled using a total of  $K(K - 1)/2$  displacement vectors where each of these vectors is defined for a particular pair of primitive objects. The Gaussian mixture model that combines multiple primitive objects is described in the next section.

### 3. GAUSSIAN MIXTURE MODEL

We model the distribution of image pixels with GMMs with  $K$  components. Each pixel can belong to one of the  $K$  Gaussian components. For each pixel  $j = 1, \dots, N$ , there is a corresponding label variable  $y_j \in \{1, \dots, K\}$ , where  $y_j = k$  denotes the event of the  $j$ 'th pixel belonging to the  $k$ 'th Gaussian component. A GMM is fully defined by the set of parameters  $\Theta = \{\alpha_1, \dots, \alpha_K, \theta_1, \dots, \theta_K\}$  where each  $\theta_k = \{\mu_k, \Sigma_k\}$  represents the parameters of the  $k$ 'th Gaussian distribution  $p_k(\mathbf{x}|\theta_k)$ .  $\mu_k \in \mathbb{R}^d$  denotes the mean vector and  $\Sigma_k \in \mathbb{S}_{++}^d$  denotes the covariance matrix of the  $k$ 'th Gaussian component. Probability of a pixel belonging to the  $k$ 'th Gaussian component is denoted by  $\alpha_k \in [0, 1]$ , where  $\alpha_1, \dots, \alpha_K$  are constrained to sum up to 1, i.e.,  $\sum_{k=1}^K \alpha_k = 1$ .

For a given set of  $N$  pixels  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  that are assumed to be independent and identically distributed (i.i.d.) according to the mixture probability density function  $p(\mathbf{x}|\Theta) = \sum_{k=1}^K \alpha_k p_k(\mathbf{x}|\theta_k)$ , the objective is to obtain the maximum likelihood estimate of  $\Theta$  by finding the parameters that maximize the log-likelihood function

$$\log L(\Theta|\mathcal{X}) = \log p(\mathcal{X}|\Theta) = \sum_{j=1}^N \log \left( \sum_{k=1}^K \alpha_k p_k(\mathbf{x}_j|\theta_k) \right). \quad (1)$$

In the compound object detection problem, we assume that we are given an example compound structure of interest. This input is expected to be in the form of delineated regions for the primitive objects. The regions corresponding to the primitive objects can be obtained using basic low-level operations such as morphological opening or closing, or can be obtained via manual selection. The total of  $\tilde{N}$  pixels belonging to the given  $K$  primitive objects are used to fit a GMM with  $K$  components to the reference structure where each primitive object is modeled by one of the Gaussian components. Since the memberships of all reference pixels to the Gaussian components are known, the reference GMM parameters can be directly obtained using maximum likelihood estimates where the parameters of the  $k$ 'th Gaussian component are computed using the pixels belonging to the  $k$ 'th region in the delineated reference structure. The resulting reference GMM is defined by its parameters  $\tilde{\Theta} = \{\tilde{\alpha}_1, \dots, \tilde{\alpha}_K, \tilde{\theta}_1, \dots, \tilde{\theta}_K\}$  where  $\tilde{\theta}_k = \{\tilde{\mu}_k, \tilde{\Sigma}_k\}$ ,  $k = 1, \dots, K$ .

A target image with  $N$  pixels is represented by  $N$   $d$ -dimensional feature vectors as described in Sec. 2. Our goal is to identify the pixels in the target image that are the most similar to the selected regions in the reference structure. This is achieved by clustering the pixels of the target image using a robust constrained GMM model with  $K$  components where  $K$  is the same as the number of components in the reference GMM. This model uses the parameters of the reference GMM to form constraints on its parameters, and selects  $\tilde{N}$  pixels as inliers while rejecting the rest as outliers. The proposed GMM model is described below. The algorithm for estimating it from image data is presented in Sec. 4.

The constraints are defined between pairs of Gaussian components, one from the reference GMM and the other one from the target GMM. In order to identify the correspondence relation between the Gaussian components of the reference and target GMMs, we define a mapping denoted as  $P: \{1, \dots, K\} \rightarrow \{1, \dots, K\}$  where  $P(i) = j$  if the  $j$ 'th Gaussian component of the reference GMM corresponds to the  $i$ 'th Gaussian component of the target GMM. Given the correspondence relation  $P$ , the constraints are defined as follows.

- We want to keep the relative sizes of the components of reference and target GMMs same, i.e.,  $\alpha_i = \tilde{\alpha}_{P(i)}$  for  $i = 1, \dots, K$ .
- We want the spectral content of the reference and target components to be similar. For this purpose, we constrain the multispectral part of each target mean to lie inside a confidence ellipsoid around the reference mean, i.e.,  $(\mu_i^{ms} - \tilde{\mu}_{P(i)}^{ms})^T (\tilde{\Sigma}_{P(i)}^{ms})^{-1} (\mu_i^{ms} - \tilde{\mu}_{P(i)}^{ms}) \leq \beta$  for  $i = 1, \dots, K$  where  $\beta$  is a constant.
- We also set the multispectral parts of the covariance matrices of the corresponding reference and target components to be the same, i.e.,  $\Sigma_i^{ms} = \tilde{\Sigma}_{P(i)}^{ms}$  for  $i = 1, \dots, K$ .
- We want to preserve the spatial layout of the reference GMM in the target GMM. Thus, first, we define  $K(K-1)/2$  displacement vectors  $\mathbf{d}_{ij}$ ,  $i = 1, \dots, K-1$ ,  $j = i+1, \dots, K$ , between the spatial parts of the reference means, i.e.,  $\tilde{\mu}_{P(i)}^{xy} + \mathbf{d}_{ij} = \tilde{\mu}_{P(j)}^{xy}$ . Then, the spatial layout of the target GMM can be expressed as  $\mu_i^{xy} + \mathbf{d}_{ij} - \mu_j^{xy} = \mathbf{t}_{ij}$  where  $\|\mathbf{t}_{ij}\|_1 \leq u$  and the constant  $u \in \mathbb{R}_+$  specifies the allowed amount of deviation from the reference spatial relations.

- Finally, we want the minimum and maximum eigenvalues of the spatial parts of the reference and target covariances to be the same so that the aspect ratio of each reference primitive object is preserved in the target, i.e.,  $\lambda_{\min}(\mathbf{\Sigma}_i^{xy}) = \lambda_{\min}(\tilde{\mathbf{\Sigma}}_{P(i)}^{xy})$  and  $\lambda_{\max}(\mathbf{\Sigma}_i^{xy}) = \lambda_{\max}(\tilde{\mathbf{\Sigma}}_{P(i)}^{xy})$  for  $k = 1, \dots, K$ .

#### 4. DETECTION ALGORITHM

The common practice for estimating the GMM parameters in the literature is to use the expectation-maximization (EM) algorithm. The EM algorithm iteratively updates the parameters of individual Gaussian distributions in the mixture. Each iteration consists of two main steps called the E-Step and the M-Step. In the E-Step, the algorithm makes a probabilistic guess for the label variables  $y_j, j = 1, \dots, N$ . These estimates, denoted as  $w_{jk}$ , are the posterior probabilities of the label variables given the corresponding data points  $\mathbf{x}_j, j = 1, \dots, N$ , i.e.,  $w_{jk} = P(y_j = k | \mathbf{x}_j, \Theta)$ . In the M-Step, the algorithm finds the parameters of the GMM using the estimated label probabilities. The Gaussian component probabilities  $\alpha_k, k = 1, \dots, K$ , are calculated as the normalized average of the posterior label probabilities. The mean vectors  $\boldsymbol{\mu}_k, k = 1, \dots, K$ , are computed as the normalized weighted averages of the data points where the weights are the label posterior probabilities. Similarly, using the same weights as in the mean calculation, the covariance matrices  $\mathbf{\Sigma}_k, k = 1, \dots, K$ , are calculated as the normalized weighted averages of the outer products of the mean subtracted data points. The EM algorithm alternates between these two steps until convergence. In general, the algorithm is run either until an allowed maximum number of iterations is attained or until the difference between the log-likelihood values at two successive iterations falls below some given threshold value. Details of the EM algorithm for GMM estimation can be found in Ref. 5.

The compound object detection algorithm uses a modified version of the standard EM algorithm to fit the robust constrained GMM described in Sec. 3 to the target image data. The input contains the reference GMM with parameters  $\tilde{\Theta} = \{\tilde{\alpha}_k, \tilde{\boldsymbol{\mu}}_k, \tilde{\mathbf{\Sigma}}_k\}_{k=1}^K$  estimated using a total of  $\tilde{N}$  pixels belonging to  $K$  primitive objects, and the target image with  $N$  pixels  $\mathbf{x}_j, j = 1, \dots, N$ . The reference spatial layout is also constructed in terms of  $K(K-1)/2$  displacement vectors  $\mathbf{d}_{ij}, i = 1, \dots, K-1, j = i+1, \dots, K$ , computed from the centroids of the primitive objects. We associate each pixel in the target image with an indicator variable  $z_j, j = 1, \dots, N$ , where  $z_j = 1$  if a pixel is determined as an inlier and  $z_j = 0$  if it is not. Then, the robust GMM model is estimated by maximizing the trimmed log-likelihood function<sup>6</sup>

$$\sum_{j=1}^N z_j \log \left( \sum_{k=1}^K \alpha_k p_k(\mathbf{x}_j | \boldsymbol{\theta}_k) \right) \quad (2)$$

with the additional constraint  $\sum_{j=1}^N z_j = \tilde{N}$ .

The proposed EM algorithm has three main differences from the standard algorithm. First, the data points  $\mathbf{x}_j, j = 1, \dots, N$ , are weighted by inlier indicator variables  $z_j, j = 1, \dots, N$ , while calculating the mean vectors and covariance matrices. Since the inlier indicator variables are either one or zero, it simply means that, the mean vectors and the covariance matrices are calculated only from the inlier data points. The second difference is that, since the parameters might not satisfy the desired constraints after being updated in the M-Step, they are projected onto feasible sets.<sup>7</sup> The third difference is the addition of the Z-Step which is used to determine the inliers. In the Z-step,  $\tilde{N}$  inlier indicator variables corresponding to the data points having the highest log-likelihoods using the parameters of the current iteration are set to one and the others are set to zero.

This procedure is run by starting from different initializations of the target GMM on the target image. For each initialization, the EM algorithm gives the GMM parameters and the indicator variables  $z_j, j = 1, \dots, N$ , corresponding to a local maximum of the trimmed likelihood function in (2). Each result corresponds to a grouping of the pixels that have high likelihoods of belonging to the reference Gaussian object models while satisfying the spatial layout constraints. The results can be sorted in descending order of the likelihood values, and a list of compound structures detected in the target image can be obtained by truncating this list at a particular likelihood value.

## 5. EXPERIMENTS

Proof-of-concept experiments were performed on an 8-band multispectral WorldView-2 image of Ankara, Turkey with  $500 \times 500$  pixels and 2 m spatial resolution. The reference compound structures were obtained by manual delineation of the individual primitive objects. The parameters of the reference Gaussian components were obtained using maximum likelihood estimation. In particular, the component probabilities ( $\tilde{\alpha}_k, k = 1, \dots, K$ ) were estimated using the ratio of the number of pixels in each primitive object to the total number of pixels in the compound structure, and the means and the covariance matrices were estimated using the pixels belonging to each primitive object. After this supervised step, the rest of the detection process was performed fully unsupervised using the EM algorithm described in Sec. 4. Note that, the algorithm does not require any initial segmentation while performing object detection because it can group individual pixels that have high likelihoods of belonging to the Gaussian object models while satisfying the spatial layout constraints.

Since each different initialization of the EM algorithm converges to a local maximum of the likelihood function and there is no additional information about the expected locations of similar compound structures in the target image, we used a straightforward initialization procedure using uniform sampling of the image coordinates. First, the reference structure was placed at the top-left corner of the target image. Then, the  $x$  and  $y$  coordinates were incremented by 25 pixels to form a grid of points that were used as offsets to be added to the centroids of the reference objects for initialization while preserving the displacement relations of the centroids computed from the reference GMM. This resulted in  $19 \times 19 = 361$  runs for the EM algorithm. For each run, after calculating the initial centroids using these offset values, the spatial covariances were initialized to the reference GMM's corresponding spatial covariances. Furthermore, the means and covariances corresponding to multispectral values were also initialized to the reference GMM's corresponding means and covariances. Similarly, the Gaussian component probabilities were initialized to reference Gaussian component probabilities. Finally, the number of inliers was set to the total number of pixels in the reference structure. For all experiments, the number of mixture components was fixed to the number of primitive objects in the reference structure.

Fig. 1 shows an example structure composed of four buildings with red roofs placed in a diamond formation. The resulting target GMMs obtained after the convergence of the EM algorithm for each of the 361 runs were ranked according to the resulting likelihood values computed as in (2). The figure shows the top nine structures that corresponded to the highest likelihood values. The spatial layout model and the constraints defined in Secs. 2 and 3, respectively, allow the individual Gaussian components to rotate around their centroids while preserving the relative displacements computed from the reference GMM. Therefore, some of the detected structures corresponded to formations by rotated buildings (e.g., cross-like formation of four buildings, and parallel groups of two buildings) where the displacements between pairwise centroids were always very similar to those in the reference structure because of the constraints used.

Fig. 2 shows another example structure corresponding to an intersection of four road segments. Similar to the previous example, the resulting target GMMs obtained after the convergence of the EM algorithm for each of the 361 runs were ranked according to the resulting likelihood values. The figure shows the top six structures that corresponded to the highest likelihood values. All results except the third one corresponded to intersections that were similar to the reference structure. The third result shows an interesting case where nearby road segments formed a different structure because of the allowed rotations around the centroids with almost identical displacement. Additional constraints can be used to restrict or relax both the appearances and the spatial layout of the primitive objects within the compound structures of interest.

## 6. CONCLUSIONS

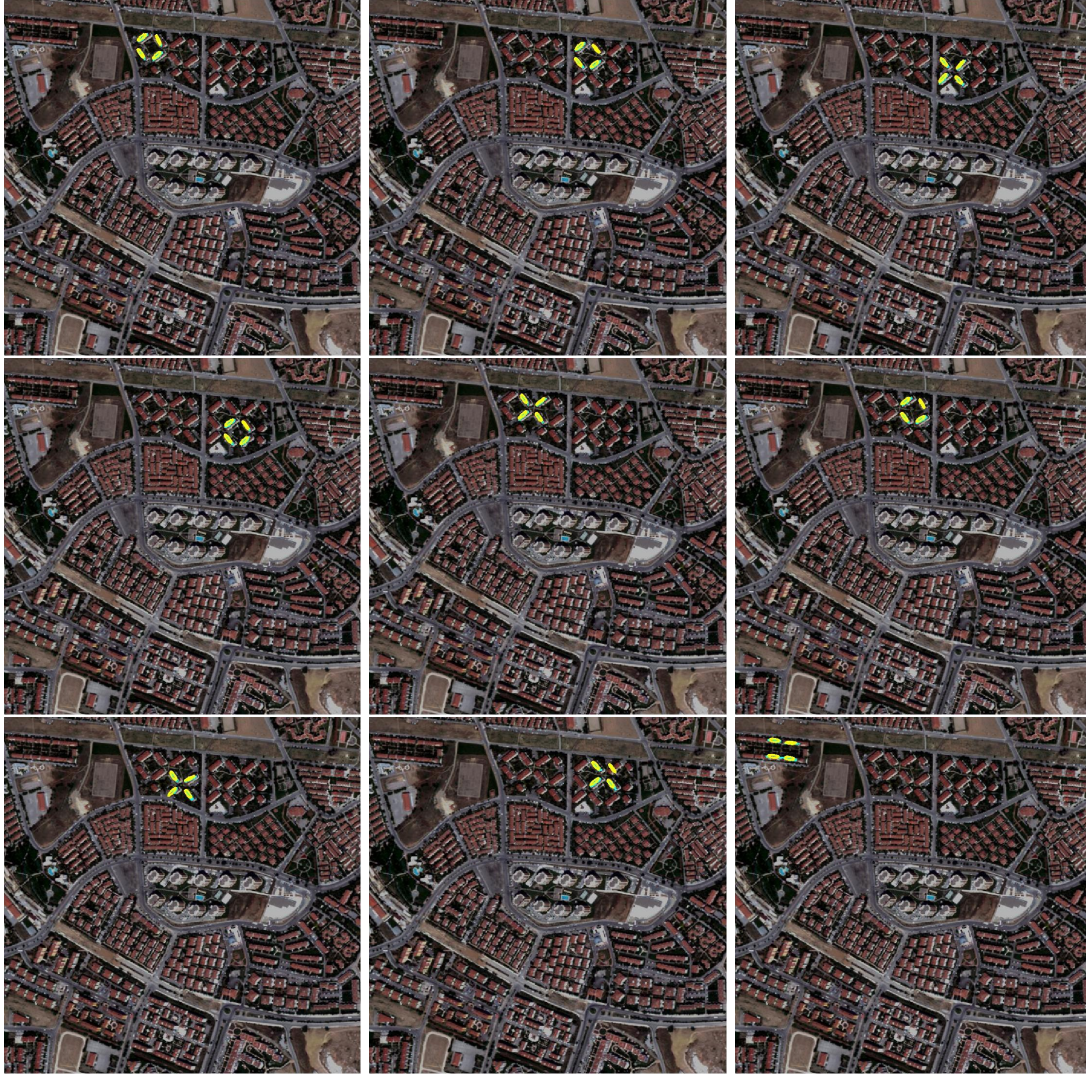
We presented a new Gaussian mixture model that uses the individual Gaussian components to represent the spectral and shape contents of basic primitive objects, and proposed a new expectation-maximization algorithm that can incorporate spectral and spatial constraints for the detection of compound structures that are comprised of spatial arrangements of such objects. Given an example compound structure of interest, first, a reference GMM was estimated from the pixels belonging to the manually delineated primitive objects. Then, the EM algorithm was used to fit a robust GMM to the target image data so that the pixels that had high likelihoods of belonging to the Gaussian object models and satisfied the spatial layout constraints could be grouped to perform unsupervised object detection.



(a) RGB image



(b) Reference structure



(c) Detected structures

Figure 1. Detection of an example structure composed of four buildings with red roofs in a diamond formation. (a) shows the RGB image formed by the visible bands. (b) shows a close up of the four patches, that were manually delineated as primitive objects, overlayed on the RGB image as yellow polygons. (c) shows the top nine structures that corresponded to the highest likelihood values at the end of all runs of the EM algorithm. For each result, the pixels selected as inliers are marked in cyan, and the resulting Gaussians are overlayed as yellow ellipses drawn at three standard deviations.

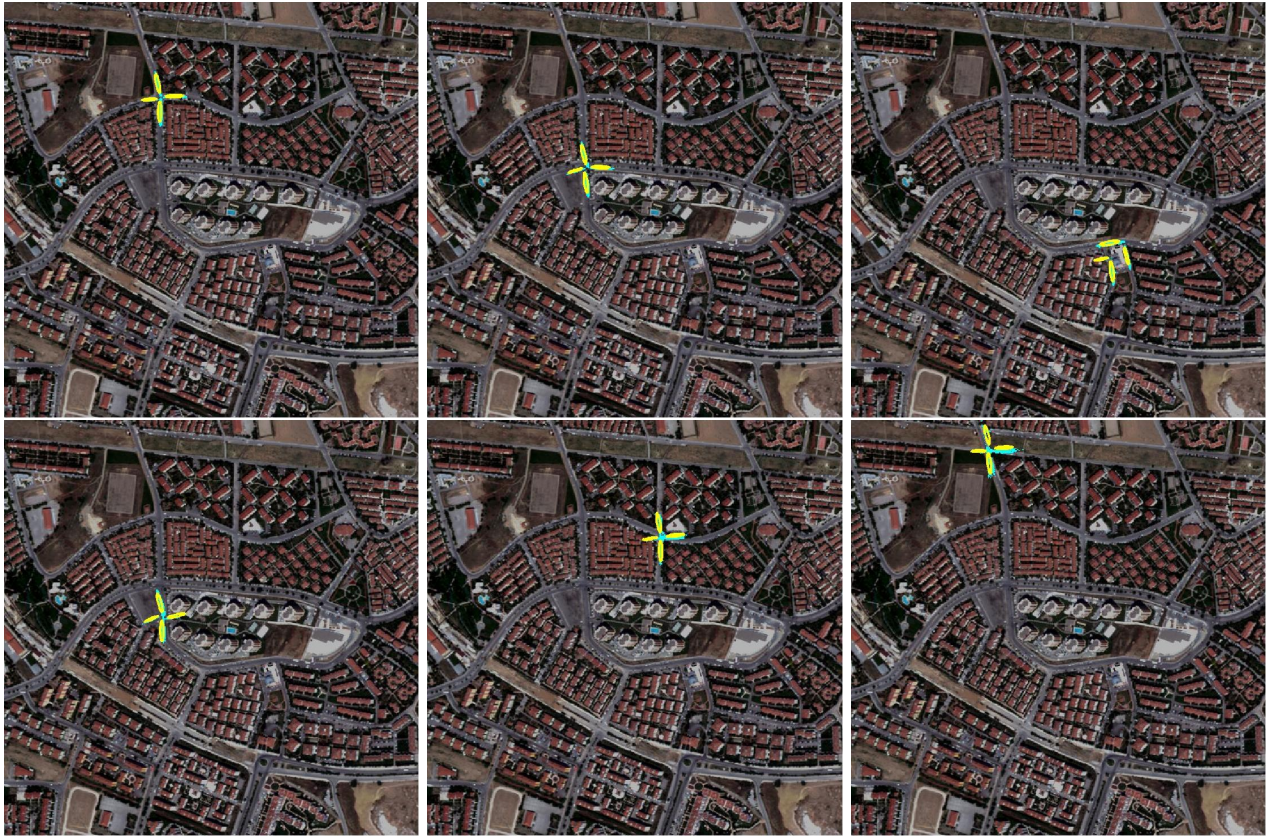




(a) RGB image



(b) Reference structure



(c) Detected structures

Figure 2. Detection of an example structure corresponding to an intersection of four road segments. (a) shows the RGB image formed by the visible bands. (b) shows a close up of the four patches, that were manually delineated as primitive objects, overlaid on the RGB image as yellow polygons. (c) shows the top six structures that corresponded to the highest likelihood values at the end of all runs of the EM algorithm. For each result, the pixels selected as inliers are marked in cyan, and the resulting Gaussians are overlaid as yellow ellipses drawn at three standard deviations.

The initial experiments showed that the proposed method can detect high-level structures that cannot be modeled using traditional techniques. Furthermore, it has a very important advantage of not requiring any initial segmentation while performing object detection by grouping individual pixels. In the proof-of-concept experiments presented in this paper, all primitive objects corresponded to the same type, i.e., buildings in Fig. 1 and roads in Fig. 2, but the algorithm can use any type of primitive object. Therefore, future work includes experiments with other types of compound structures in larger data sets. We are also planning to extend the model with additional constraints.

## ACKNOWLEDGMENTS

This work was supported in part by the TUBITAK Grant 109E193.

## REFERENCES

- [1] Gaetano, R., Scarpa, G., and Poggi, G., “Hierarchical texture-based segmentation of multiresolution remote-sensing images,” *IEEE Transactions on Geoscience and Remote Sensing* **47**, 2129–2141 (July 2009).
- [2] Zamalieva, D., Aksoy, S., and Tilton, J. C., “Finding compound structures in images using image segmentation and graph-based knowledge discovery,” in [*Proceedings of IEEE International Geoscience and Remote Sensing Symposium*], **V**, 252–255 (July 13–17, 2009).
- [3] Vanegas, M. C., Bloch, I., and Inglada, J., “Detection of aligned objects for high resolution image understanding,” in [*Proceedings of IEEE International Geoscience and Remote Sensing Symposium*], 464–467 (July 25–30, 2010).
- [4] Akcay, H. G. and Aksoy, S., “Detection of compound structures using hierarchical clustering of statistical and structural features,” in [*Proceedings of IEEE International Geoscience and Remote Sensing Symposium*], 2385–2388 (July 25–29, 2011).
- [5] McLachlan, G. and Peel, D., [*Finite Mixture Models*], John Wiley & Sons, Inc. (2000).
- [6] Hadi, A. S. and Luceno, A., “Maximum trimmed likelihood estimators: A unified approach, examples, and algorithms,” *Computational Statistics & Data Analysis* **25**, 251–272 (August 1997).
- [7] Bertsekas, D. P., [*Nonlinear Programming*], Athena Scientific (1999).