



Automatic segmentation of colon glands using object-graphs

Cigdem Gunduz-Demir^{a,*}, Melih Kandemir^a, Akif Burak Tosun^a, Cenk Sokmensuer^b

^a Department of Computer Engineering, Bilkent University, Ankara TR-06800, Turkey

^b Department of Pathology, Hacettepe University Medical School, Ankara TR-06100, Turkey

ARTICLE INFO

Article history:

Received 28 August 2008

Received in revised form 24 July 2009

Accepted 10 September 2009

Available online 19 September 2009

Keywords:

Gland segmentation

Image segmentation

Histopathological image analysis

Object-graphs

Attributed graphs

Colon adenocarcinoma

ABSTRACT

Gland segmentation is an important step to automate the analysis of biopsies that contain glandular structures. However, this remains a challenging problem as the variation in staining, fixation, and sectioning procedures lead to a considerable amount of artifacts and variances in tissue sections, which may result in huge variances in gland appearances. In this work, we report a new approach for gland segmentation. This approach decomposes the tissue image into a set of primitive objects and segments glands making use of the organizational properties of these objects, which are quantified with the definition of object-graphs. As opposed to the previous literature, the proposed approach employs the object-based information for the gland segmentation problem, instead of using the pixel-based information alone. Working with the images of colon tissues, our experiments demonstrate that the proposed object-graph approach yields high segmentation accuracies for the training and test sets and significantly improves the segmentation performance of its pixel-based counterparts. The experiments also show that the object-based structure of the proposed approach provides more tolerance to artifacts and variances in tissues.

© 2009 Elsevier B.V. All rights reserved.

1. Introduction

Histopathological examination includes examining a biopsy tissue under a microscope for the identification of tissue changes associated with disease. In the current practice of medicine, this examination is the most important tool for routine clinical diagnosis of a large group of diseases including cancer. However, as it mainly relies on the visual interpretation of a pathologist, it may lead to a certain level of subjectivity (Thomas et al., 1983; Andrión et al., 1995). To help pathologists in diagnosis, and hence, to reduce the subjectivity level, it has been proposed to use computational methods that provide objective measures (Wolberg et al., 1995; Thiran and Macq, 1996; Choi et al., 1997; Hamilton et al., 1997; Esgiar et al., 1998; Spyridonos et al., 2001; Wiltgen et al., 2003; Nielsen et al., 1999; Esgiar et al., 2002; Weyn et al., 1999; Keenan et al., 2000; Demir et al., 2005; Gunduz-Demir, 2007). These computational methods extract a set of mathematical features (e.g., morphological (Wolberg et al., 1995; Thiran and Macq, 1996), textural (Hamilton et al., 1997; Choi et al., 1997; Esgiar et al., 1998; Spyridonos et al., 2001; Wiltgen et al., 2003), fractal (Nielsen et al., 1999; Esgiar et al., 2002), and structural (Choi et al., 1997; Weyn et al., 1999; Keenan et al., 2000; Demir et al.,

2005; Gunduz-Demir, 2007)) from a tissue image for its quantification and use these mathematical features to objectively measure the degree of the tissue changes associated with a disease of the interest. Different types of features might be necessary to quantify the tissue changes as these changes show differences from one tissue type to another as well as from one disease type to another. For example, soft tissue tumors change the cell distribution in the tissue whereas adenocarcinomas change the architecture of glands¹ as well. To identify the latter type of neoplastic diseases, which cause changes in gland architectures, the very first step is to segment the tissue into its gland structures.

In literature, there are few studies that focus on the problem of automatic gland segmentation for tissues that contain gland structures (Wu et al., 2005a,b; Naik et al., 2007; Farjam et al., 2007). These studies make use of the fact that glands are characterized by their luminal areas surrounded by epithelial cells; an example of the histopathological image of a colon tissue is given in Fig. 1. In order to capture this characterization, these studies first identify the pixels of different classes (e.g., nucleus, stroma, and lumen classes) and then form gland regions using this class information of pixels. For example, in Wu et al. (2005a), nucleus pixels are identified applying a threshold to the intensities of pixels after they are convolved with a composition of directional filters. The regions

* Corresponding author. Tel.: +90 312 290 3443; fax: +90 312 266 4047.

E-mail addresses: gunduz@cs.bilkent.edu.tr (C. Gunduz-Demir), melih@cs.bilkent.edu.tr (M. Kandemir), tosun@cs.bilkent.edu.tr (A.B. Tosun), csokmensuer@hacettepe.edu.tr (C. Sokmensuer).

¹ Many types of tissues such as colon, prostate, breast, and lung include glands. Neoplastic diseases that originate from these tissues cause structural and organizational changes in their glands.

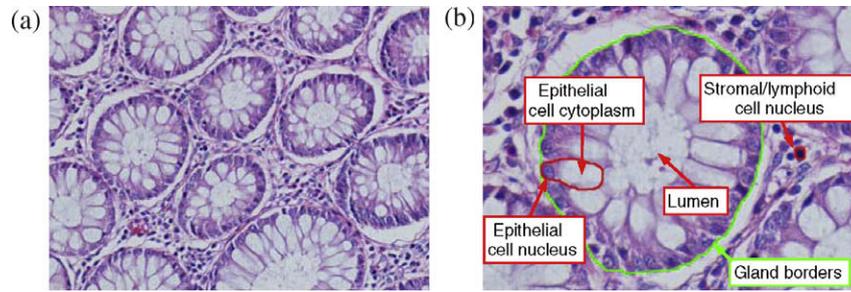


Fig. 1. (a) A histopathological image of a colon tissue, which is stained with the routinely used hematoxylin-and-eosin technique, and (b) an individual gland of a colon tissue.

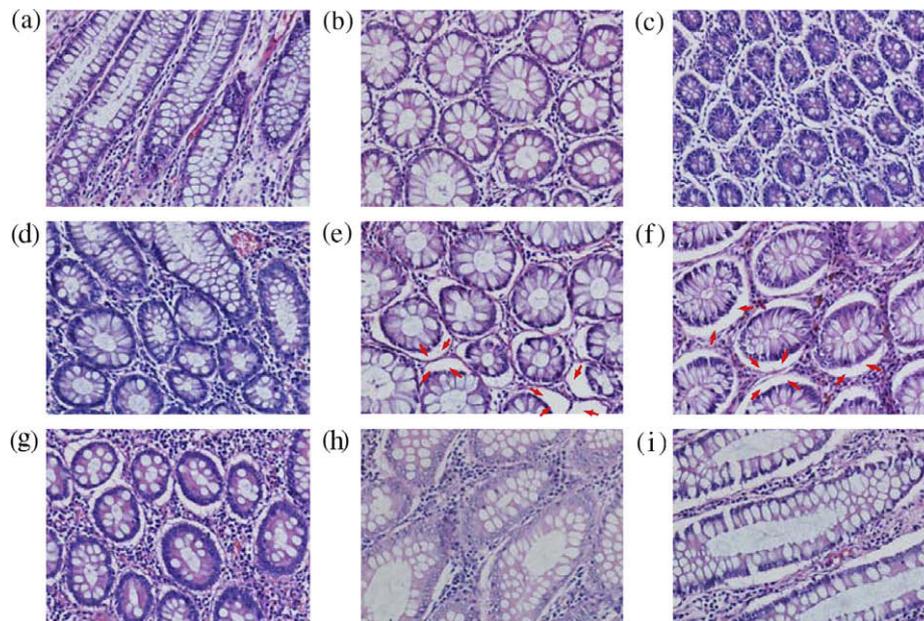


Fig. 2. Histopathological images of colon tissues, which are stained with the routinely used hematoxylin-and-eosin technique. All of the images are taken with the same magnification and the same lightning conditions.

surrounded by these pixels are then determined as glands, provided that their areas are larger than a threshold. In another work of the same authors (Wu et al., 2005b), nucleus and lumen pixels are first determined also applying a threshold to their intensities. Subsequently, large enough connected components of lumen pixels are identified as gland seeds and these seeds are then iteratively grown until a barrier of nuclei chain is reached. In another work (Naik et al., 2007), a Bayesian classifier is used to classify the pixels into nucleus, lumen, and cytoplasm classes based on their intensity values. Then candidate gland regions are defined as the connected components of pixels for which the classifier outputs posteriors greater than a threshold for the lumen class. Finally, false glands are eliminated according to their sizes and the probability of their surrounding pixels belonging to the cytoplasm class. In Farjam et al. (2007), after clustering the pixels into nucleus, stroma, and lumen classes based on their textural properties, the glands are obtained excluding the regions containing nucleus pixels from those containing stroma and lumen pixels.

These studies yield promising results for especially tissues in which the glands appear in more regular structures showing less variations. However, due to staining, fixation, and sectioning procedures, there is a considerable amount of artifacts and variances in tissue sections, which may result in huge variances in gland appearances. First, glands could be of different sizes, depending on the orientation of the tissue at the time of sectioning. For exam-

ple, although they are taken with the same magnification, the images shown in Fig. 2a–c have glands of different sizes. Furthermore, the improper orientation of the tissue produces tangential sectioning, which results in glands of different sizes within the same tissue image (Fig. 2d). Therefore, in false gland elimination, it is almost impossible to find an area threshold that applies for all images. Second, because of the density difference between the glandular and connective tissue structures, the fixation and sectioning procedures may result in large white artifacts on the boundaries of the glands (some of these artifacts are shown with red² arrows in Fig. 2e and f). Considering only the pixel-based information, it is more difficult to distinguish such white artifacts from luminal regions. Third, the thickness of a tissue section and the freshness of dye cause variations in the intensity distribution of a tissue image. Moreover, stain fades in time. Therefore, a single threshold value could not be found for all images to determine their nucleus pixels. Even such a threshold is manually selected or automatically determined (e.g., by the Otsu method (Otsu, 1979)) for each image, the resulting nucleus pixels do not usually form a closed component even after postprocessing the pixels (e.g., using mathematical morphology (Serra, 1982)). Thus, it is rare to find continuous

² For interpretation of color in Figs. 1–9, the reader is referred to the web version of this article.

nucleus pixels that surround the luminal area. For example, although it is more possible to find such nucleus pixels in the tissue shown in Fig. 2g, it is much more difficult to find them for tissues shown in Fig. 2h and i. Because of all these issues, using only the pixel-based information leads to incorrect gland segmentations for especially tissues with artifacts and variations.

In this paper, we report a new gland segmentation algorithm that relies on decomposing the image into a set of primitive objects (nucleus and lumen objects) and then making use of the organizational properties of these objects instead of using the pixel-based information alone. This object-based algorithm is a three-step region growing approach. First, it constructs a graph on all of its objects and determines gland seeds based on the features extracted from this object-graph. Then, it constructs another graph, this time on its nucleus objects, and uses this second object-graph for growing the gland seeds. Finally, it determines the final boundary of glands based on the locations of the nucleus objects. After this region-growing process, false glands are eliminated based on the cluster information of the grown regions. Working with colon tissues of 36 different patients, our experiments show that the region-growing process leads to 82.57% average segmentation accuracy on the test set and that this accuracy increases to 87.59% after the false gland elimination step. These results (both before and after false gland elimination) demonstrate that the proposed object-based algorithm significantly improves the segmentation performance of its pixel-based counterparts. To the best of our knowledge, this is the first demonstration of the use of object-graphs for the purpose of gland segmentation.

2. Object-graph approach

2.1. Overview

The proposed object-graph approach relies on modeling the regular structure of glands. For this purpose, it decomposes a tissue image into a set of objects, which represent different tissue components, and uses the way that they distribute within the tissue in a region-growing process to determine the locations of gland structures. Compared to pixel-based information, the use of object-based information in a region-growing process yields more robust segmentations as pixel intensities are expected to be more sensitive to the noise that arises from the staining, fixation, and sectioning related problems.

In Fig. 3, a schematic of the proposed approach is provided. In this approach, an image is first decomposed into its tissue components. Since it is very difficult to exactly locate the components, they are approximately represented transforming the image into a set of circular objects (nucleus and lumen objects). For the image given in Fig. 1a, these circular objects are shown in the first step of Fig. 3; here the nucleus and lumen objects are shown with black and cyan, respectively.

After the transformation, gland segmentation is achieved by making use of the organizational properties of these objects, which are quantified with the definition of object-graphs. This gland segmentation algorithm includes a three-step region growing process (with the initial gland seed determination, gland seed growing, and gland boundary detection steps) followed by false gland elimination. In initial gland seed determination, the lumen objects are divided into two classes based on their organizational properties: the “gland” class corresponding to the lumen objects inside a glandular region and the “non-gland” class corresponding to those outside a glandular region. The lumen objects falling in the gland class are identified as initial gland seeds. In the second step of Fig. 3, the

lumen objects that are considered as initial gland seeds are shown in red whereas the other lumen objects are shown in green.

The initial gland seeds are grown to identify inner gland regions. In region growing, one has to determine the locations where the growing process is supposed to stop. The object-graph approach uses the nucleus objects to find these locations since the inner gland regions are expected to be surrounded by nuclei. To this end, it constructs another object-graph on the nucleus objects and uses the edges of this graph (i.e., the pixels that correspond to these edges) to stop the region growing process. The graph edges used to stop region growing and the inner regions obtained with this process are illustrated in the third step of Fig. 3. Subsequently, the outer boundary is found by extending the inner region of a gland to include the nucleus objects that are in its close proximity. The outer gland boundaries are shown in the fourth step of Fig. 3.

After growing the seeds, the ones that do not show the characteristics of a colon gland are eliminated. For a colon gland, its inner part is expected to contain the luminal region of the gland and epithelial cell cytoplasm whereas its outer part is expected to contain epithelial cell nuclei. Thus, the proposed algorithm divides an identified gland region into its inner and outer parts and extracts a set of features using their cluster information. Using these features, false glands are eliminated in a supervised manner. The final gland locations obtained after false gland elimination are shown in the fifth step of Fig. 3. In the next subsections, the details of these steps are explained.

2.2. Object definition

The proposed approach defines objects to represent tissue components. In the definition of the objects, the ideal way would be to identify different components (such as epithelial cell nuclei, epithelial cell cytoplasm, stromal cell nuclei, and lumina) in the tissue. However, this would require segmenting these tissue components, which gives rise to more difficult segmentation problem. Therefore, instead of exactly determining their locations, we approximately represent these components by transforming the image into a set of circular primitives.³ In this work, two different types of circular objects are defined: one for representing cell nuclei and the other for representing lumina and epithelial cell cytoplasm (herein referred to as “nucleus objects” and “lumen objects”, respectively). For defining the circular objects, the “circle-fit algorithm” that we propose in our previous work (Tosun et al., 2009) has been employed.

2.2.1. Circle-fit algorithm

On a given set of pixels $\mathcal{P} = \{x_i\}$, the circle-fit algorithm locates a set of circles iteratively. In the first step of this algorithm, each pixel x_i is assigned to the largest possible circle that includes only the pixels of \mathcal{P} and the pixel x_i . In its second step, the pixels assigned to the same circle are connected and the connected components smaller than an area threshold are eliminated. Then, for each connected component with a size greater than the threshold, the first two steps are iteratively repeated (considering only the pixels of this component) until there is no change in the subsequent iterations. Note that there will be no change, if the component is circular.

In this work, the circle-fit algorithm is run twice for the given pixels \mathcal{P} . First, it is run on all pixels in \mathcal{P} to find a set of circles. Then, it is run again on the pixels that are in \mathcal{P} but not belong to any circles found in the first run. Finally, the circles found in the first and second runs are merged. In Fig. 4, the result of the

³ Here a circular shape is particularly selected for the transformation as the borders of tissue components typically comprise curves and circular shapes are found more efficiently compared to for example elliptical shapes.

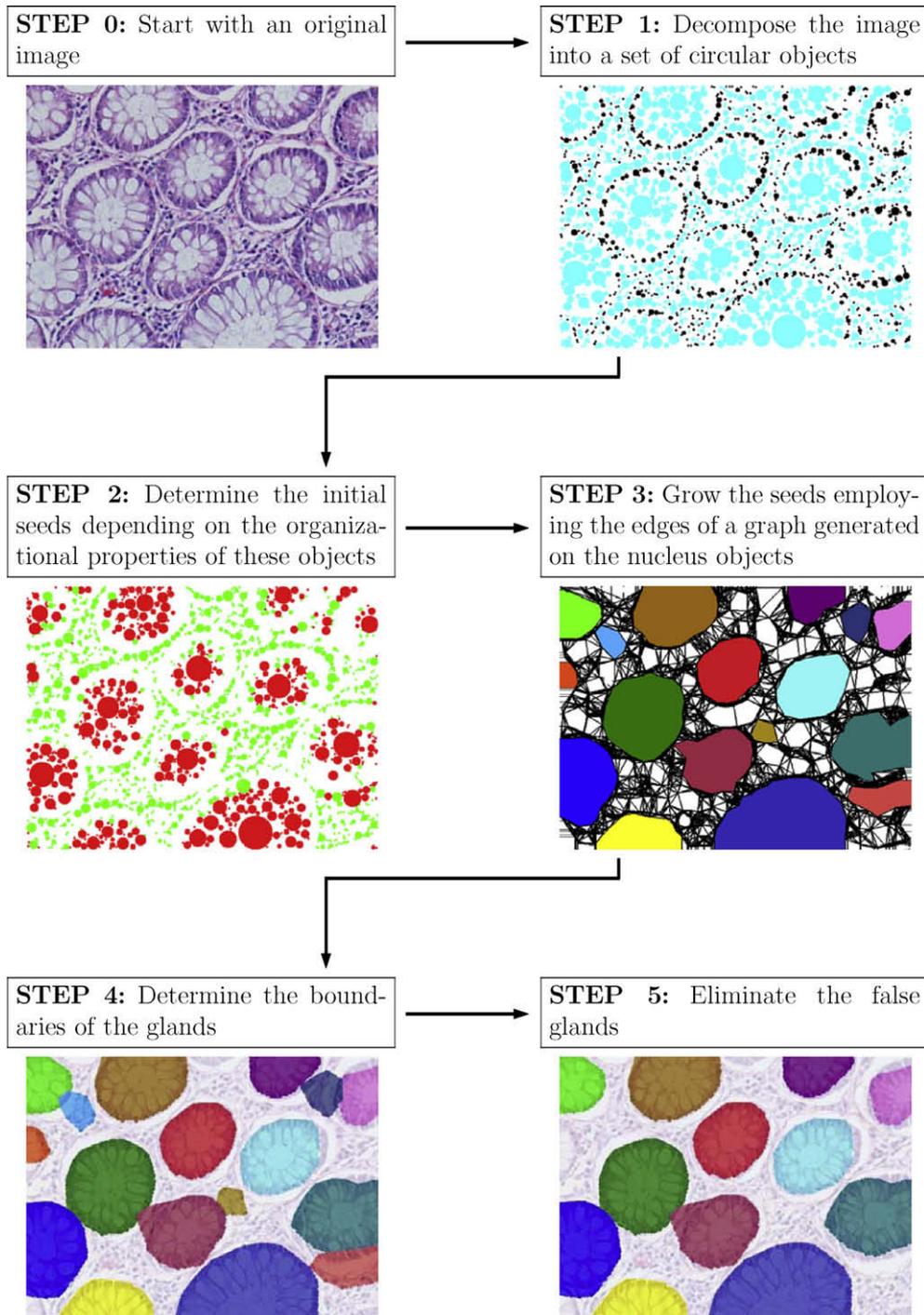


Fig. 3. Overview scheme of the proposed method.

circle-fit algorithm is illustrated on a small image; in this figure, the circles shown with red and yellow correspond to the circles found in the first and second runs, respectively.

2.2.2. Image decomposition

For a given tissue image, the pixels are first quantized into three clusters using the k -means algorithm. In this work, the number of clusters is particularly selected as three since there are mainly three color groups in the image of a tissue stained with hematoxylin-and-eosin. These colors are purple that correspond to nucleus pixels, pink that corresponds to stroma pixels, and white that corresponds to lumen and epithelial cell cytoplasm pixels.

After color quantization, the circle-fit algorithm is run for the nucleus and lumen⁴ clusters, separately. Before calling the circle-fit algorithm, morphological operators are applied to the pixels of each cluster to reduce the noise that arises from the incorrect assignment of pixels in color quantization. This approach only considers the nucleus and lumen clusters but not the stroma cluster since the experiments demonstrate that the consideration of the stroma cluster does not improve the performance of the system. Thus, the stroma cluster is not considered for the sake of simplicity.

⁴ For the sake of simplicity, we refer the cluster that corresponds to lumen and epithelial cell cytoplasm pixels as “lumen” cluster.

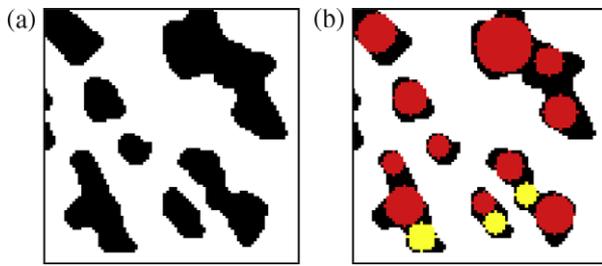


Fig. 4. The result of the circle-fit algorithm: (a) pixels in the given set are shown with black and (b) circles found in the first and second runs are shown with red and yellow, respectively. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

In our previous work (Tosun et al., 2009), such a transformation is also used for segmenting a low-magnification tissue image into its homogeneous regions, which consist of either cancerous or normal parts. For example, in Fig. 5a, one of such images and its segmentation result obtained by the algorithm that is proposed in Tosun et al. (2009) are shown.⁵ This previous work introduces a new homogeneity measure that quantifies how uniform the circles are distributed in size and in space. For each particular pixel, this homogeneity measure is calculated over a window centered at this particular pixel and then segmented regions are formed by connecting pixels based on their homogeneity values. As opposed to our previous work, this paper proposes a scheme for segmenting a higher-magnification tissue image into its gland structures. In this newly proposed scheme, object-graphs are constructed on the circles of nucleus and lumen clusters, for the first time, and these graphs are used for gland segmentation. In Fig. 5b, one of such images and its segmentation result obtained by the algorithm that is proposed in this work are shown.

2.3. Initial gland seed determination

An object-graph is constructed and the local features extracted from this graph are used to determine the initial seed locations. In the object-graph construction, nucleus and lumen objects are defined as nodes and edges are assigned between each lumen object and its N -closest lumen and N -closest nucleus objects. For each lumen object L_i , the following set of local object-graph features is extracted, as illustrated in Fig. 6.

- The area of L_i .
- The areas of the lumen neighbors of L_i .
- The areas of the nucleus neighbors of L_i .
- The lengths of the edges between L_i and its lumen neighbors.
- The lengths of the edges between L_i and its nucleus neighbors.
- The angles between the edges of the lumen neighbors of L_i .
- The angles between the edges of the nucleus neighbors of L_i .

These extracted local features are used by the k -means algorithm to quantize the lumen objects into two clusters. These two clusters are automatically associated with the “gland” and “non-gland” classes using the observation that the lumen objects of the gland class are usually larger than those of the non-gland class. Thus, the average area for the first and the second cluster is computed and the cluster that has the larger average area is associated with the gland class and the remaining cluster is associated with the non-gland class. Finally, the lumen circles classified with the non-gland class are eliminated and the remaining lumen circles,

which are classified with the gland class, are identified as initial gland seeds.

The intuition behind the use of local object-graph features for the lumen object classification is that the relative spatial distribution of its closest objects to a lumen differs inside and outside the glandular regions. For example, for a lumen object inside the glandular region, its closest nucleus objects, which are expected to correspond to epithelial cells, generally locate on one side of the lumen object. On the other hand, for a lumen object outside the glandular region, its closest nucleus objects, which are expected to correspond to both epithelial and stromal cells, generally spread homogeneously around the lumen object.

The graph constructed in this step is an example of an attributed graph, in which the nodes represent image primitives and the edges represent the relations between these primitives. Attributed graphs have been shown to be effective in representing structural knowledge and have been used in many computer vision applications including structural matching (Christmas et al., 1995), similarity searching (Petraakis and Faloutsos, 1997), object recognition (Sanfeliu et al., 2002; Ahmadyfard and Kittler, 2003), object tracking (Tang and Tao, 2008), and face recognition (Luo et al., 2006). These applications mainly rely on comparing the attributed graphs of different structures with a graph matching algorithm, which has high computational complexity (Jain and Wysotzki, 2004). An object-graph can be considered as an attributed graph in the sense that its nodes correspond to circular objects (nucleus and lumen objects) and its edges represent N closeness relation between these objects. Nevertheless, the task of clustering the lumen objects into the gland and non-gland classes does not require any matching algorithms, and thus, it does not suffer from the high computational complexity of these algorithms.

2.4. Gland seed growing

In order to find the inner regions of glands, the initial gland seeds are grown employing another object-graph. This object-graph is constructed considering the nucleus objects as nodes and assigning edges between each node and its M -closest nodes.⁶ Starting from the initial seeds, regions are then grown until a graph edge is encountered (i.e., until a pixel that is located on an edge is found). The proposed approach uses the nucleus object-graph edges rather than nucleus pixels to stop region growing. This is because of the fact that nucleus pixels that surround a gland do not always form a closed component and there could be gaps between these pixels. Furthermore, due to the sectioning, fixation, and staining related problems, the size of these gaps could be very large for some images. The elimination of such gaps results in eliminating the inner regions of smaller glands as well, and hence, these gaps could not be eliminated by using the same method for all images (e.g., a morphological dilation operator with the same structuring element). In contrast, assigning edges between nucleus objects allows to define reasonable barriers for region growing for even such images.

At the end of the region growing process, there may exist very small regions. These regions are typically grown from small and isolated lumen objects, which are incorrectly classified with the gland class in the previous step. Such regions are eliminated applying a threshold to their areas. Since different images have glands of different sizes, this threshold is selected as a function of the largest region in the image (as the P percentage of the area of the largest

⁵ In Fig. 5a, the right region corresponds to the normal whereas the other two regions correspond to the cancerous.

⁶ The growing process may cause some flooding problems for boundary glands since an image usually includes only a part of a nucleus at the image boundaries and it is not always possible to define an object for such a nucleus. To avoid this problem, for a graph node, four virtual nucleus objects are defined at the image boundaries (left, right, top, and bottom boundary points, respectively) and these virtual objects are considered in the selection of the M -closest neighbors of the node.

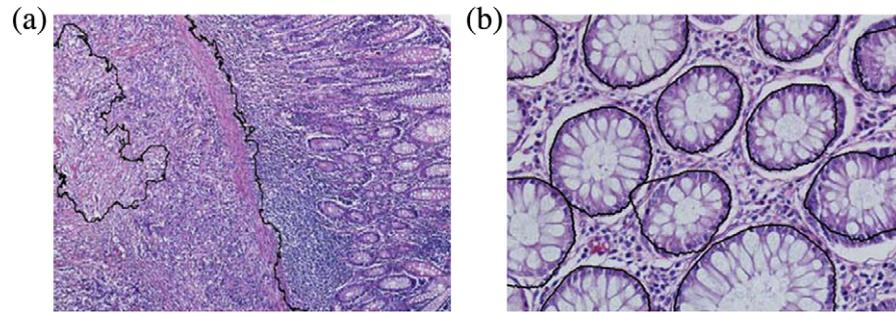
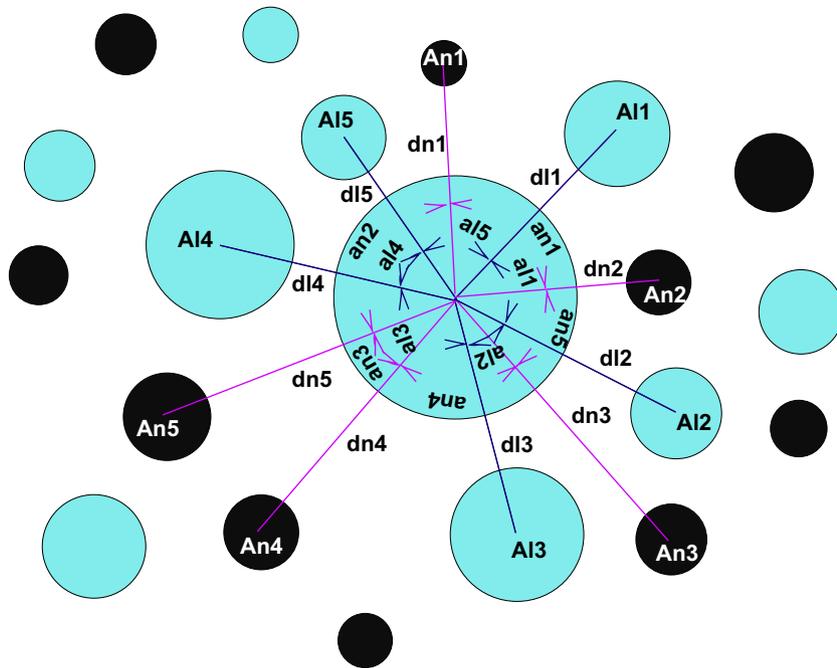
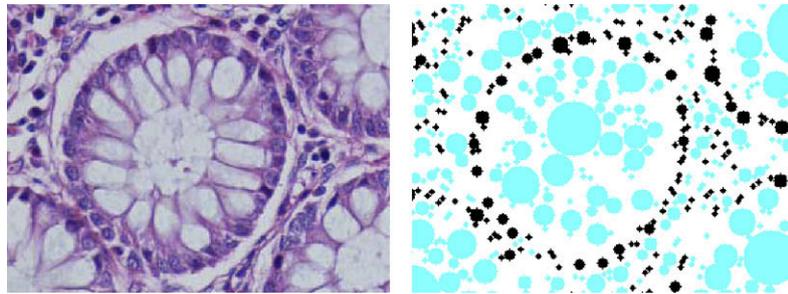


Fig. 5. (a) The segmentation of a low-magnification tissue image that is obtained by the algorithm that is proposed in our previous work (Tosun et al., 2009) and (b) the gland segmentation of a higher-magnification image that is obtained by the algorithm that is proposed in this work.



- An1 – An5: Areas of the nucleus objects
 AI1 – AI5 : Areas of the lumen objects
 dn1 – dn5 : Lengths of the edges between the lumen and its nucleus neighbors
 dl1 – dl5 : Lengths of the edges between the lumen and its lumen neighbors
 an1 – an5 : Angles between the edges of the nucleus neighbors
 al1 – al5 : Angles between the edges of the lumen neighbors

Fig. 6. The illustration of local object-graph features for a single lumen object when N is selected to be 5.

region). Furthermore, since there could be glands of different sizes within the same image, it is preferred to select smaller values of P . Note that false glands with larger sizes are to be eliminated in the false gland elimination step.

2.5. Gland boundary detection

The gland boundaries are determined by including their nuclei to the inner regions. For a gland, these nuclei are the objects that

are in a close proximity with the inner region of this gland. To find these objects, the inner region is dilated (by an amount in which the largest nucleus object in the image could be located) and nucleus objects any pixel of which is found in this dilated region are considered.

After identifying the nucleus objects, their centroids are sorted with respect to their polar angles that they make with the inner region centroid. This gives an ordered set of points (nucleus object centroids), and hence, a simple polygon; i.e., the polygon which is formed by connecting the centroids in the specified order. This simple polygon may consist of some undesired concavities. To eliminate such concavities, and thus, to obtain more accurate gland borders, this polygon is simplified by connecting each centroid to its k preceding and k succeeding centroids. As polygons are obtained connecting the centroids of nucleus objects but not their exact borders, the obtained regions are expected not to contain the half of these nucleus objects. Therefore, each polygon is dilated with the half of the largest nucleus circle in the image.

2.6. False gland elimination

Regions that do not correspond to true glands are eliminated. For this purpose, a set of features is extracted to characterize the regions and a decision tree classifier is trained in a supervised manner. For feature extraction, each gland region is divided into two: the outer part corresponding to epithelial cell nuclei of the gland and the inner part corresponding to epithelial cell cytoplasm and lumina of the gland. The width of the outer region is selected such that the largest nucleus object could fit in this outer region. Then, using the cluster information that is obtained by the k -means algorithm in the image decomposition step, the following set of features is extracted:

- The area of the outer region.
- The percentage of the nucleus cluster in the outer region.
- The percentage of the stroma cluster in the outer region.
- The percentage of the lumen cluster in the outer region.
- The area of the inner region.
- The percentage of the nucleus cluster in the inner region.
- The percentage of the stroma cluster in the inner region.
- The percentage of the lumen cluster in the inner region.

After extracting these features, each gland region in the training set is labeled as “true-gland” or “false-gland”. Since it is necessary to label gland regions lots of times (for parameter analysis), the labeling process is semi-automated using the gold standard provided by our MD collaborator. To this end, the centroid of each gland region is found and it is labeled with the true-gland class if its centroid belongs to a true gland region in the gold standard. Otherwise, it is labeled with the false-gland class. Once the decision tree classifier is trained on the training images, the rules generated by this classifier are used for false gland elimination of the other images to obtain their final gland locations.

3. Experiments

The experiments are conducted on 72 microscopic images of colon biopsy samples of 36 randomly chosen patients (two randomly selected images for each patient) from the Pathology Department archives in Hacettepe University School of Medicine. Each sample consists of 5 μm -thick tissue section and is stained with the hematoxylin-and-eosin technique. The images of these samples are taken using a Nikon Coolscope Digital Microscope with 20 \times microscope objective lens. These images are taken in the RGB color space and then converted to the Lab color space for further processing. The image resolution is 480 \times 640.

In false gland elimination, a decision tree classifier is trained to learn the rules for eliminating false glands. For this purpose, the images are divided into the training and test sets. The training set consists of 24 images of 12 patients and the test set consists of 48 images of the remaining 24 patients. The samples in the training set are used to train the decision tree classifier; the test samples are not used in training at all.

In the proposed algorithm, there are five free model parameters: (i) the area threshold in the circle-fit algorithm, (ii) N number of the closest circles of a lumen object in initial gland seed determination, (iii) M number of the closest circles of a nucleus object in nuclei graph construction, (iv) P threshold percentage to eliminate small areas in gland seed growing, and (v) k number of the connections between adjacent nucleus objects for polygon simplification. In the experiments, these parameters are selected as follows: the area threshold is 10 pixels, the number of lumen neighbors N is 5, the number of nuclei neighbors M is 10, the small

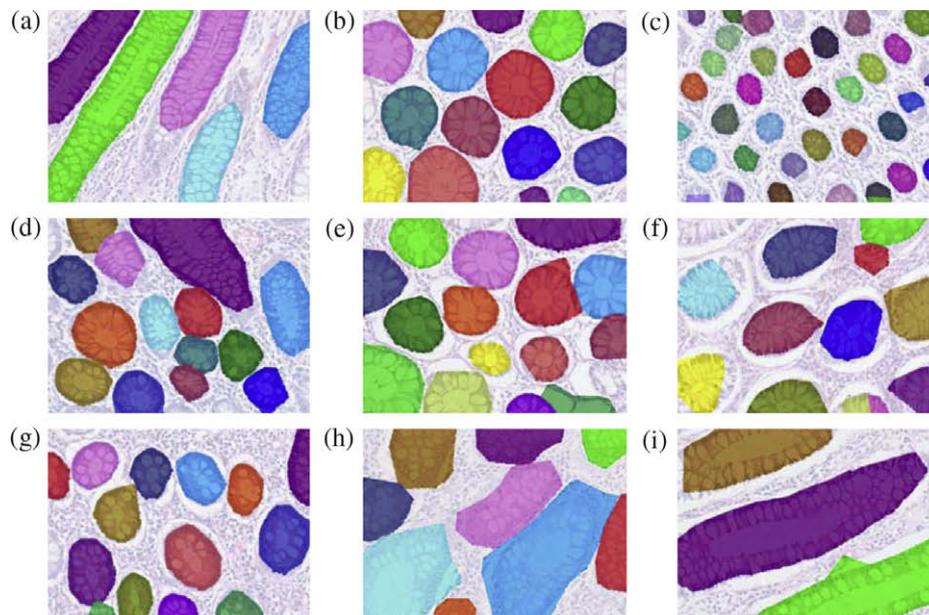


Fig. 7. The visual results obtained by the object-graph approach for the tissue images given in Fig. 2.

Table 1

For the object-graph approach, the average and the standard deviation of the sensitivity, specificity, accuracy, and Dice similarity index percentages. These results are obtained *after* applying a decision tree classifier to eliminate false glands.

	Sensitivity	Specificity	Accuracy	Dice index
Training set	83.43 ± 7.73	92.30 ± 5.78	88.00 ± 4.16	88.46 ± 4.62
Test set	85.80 ± 6.71	89.14 ± 10.40	87.59 ± 5.01	88.91 ± 4.63

Table 2

For the object-graph approach, the average and the standard deviation of the sensitivity, specificity, accuracy, and Dice similarity index percentages. These results are obtained *before* applying a decision tree classifier to eliminate false glands.

	Sensitivity	Specificity	Accuracy	Dice index
Training set	89.61 ± 4.28	65.09 ± 25.44	77.39 ± 14.99	81.73 ± 12.71
Test set	90.62 ± 5.44	72.80 ± 15.38	82.57 ± 8.36	85.59 ± 7.73

object threshold P is 5%, and the simplification factor k is 5. The selection of these parameters and their effects to the segmentation results are further discussed in the next subsection.

In Fig. 7, the segmentation results of the proposed algorithm (object-graph approach) are visually illustrated for the images given in Fig. 2. These results demonstrate that the proposed object-based algorithm leads to good segmentation results for all these tissue images even though these images have high variations and their glands appear in less regular structures. To quantitatively measure the success of these segmentation results, the true positive, false positive, true negative, and false negative rates are calculated using the manual segmentation as the gold standard⁷ and then the sensitivity, specificity, accuracy, and Dice similarity index values are computed for each of the 72 images. In Table 1, the average and standard deviation of these values are reported for the training and test set images. This table shows that the proposed algorithm yields high accuracies of 88.00% and 87.59% for the training and test sets, respectively. Moreover, it leads to Dice similarity indices greater than 88% and sensitivity and specificity rates greater than 83%.

To investigate the effects of the false gland elimination step, the proposed algorithm is also run without applying a decision tree classifier to eliminate false glands. In Table 2, the average sensitivity, specificity, accuracy, and Dice similarity index percentages that are obtained in these runs are given for the training and test sets. Here the results are provided for the training and test sets separately for better comparison although there is no training involved before applying the decision tree classifier. This table shows that without having false gland elimination, the sensitivity increases whereas the specificity, accuracy, and Dice similarity index decrease. Note that the number of positive pixels decreases after false gland elimination because some of the identified glands are eliminated. This decreases the true positive rate, and hence, the sensitivity. The Wilcoxon test with a significance level of 0.05 shows that the results before and after false gland elimination are statistically significant. The results in Table 2 demonstrate that false gland elimination is one of the important steps of the proposed gland segmentation algorithm, as in the case of previous approaches. For example, in Wu et al. (2005b), after gland segmentation, false

intestinal glands are eliminated if their surrounding nucleus pixels are not wide enough. Similarly, in Naik et al. (2007), false prostate glands are eliminated according to their sizes and the probability of their surrounding pixels belonging to the cytoplasm class.

3.1. Parameter analysis

Next, the effects of each parameter to the segmentation performance are investigated. For that, four of the five parameters are fixed and the sensitivity, specificity, accuracy, and Dice similarity index percentages are observed as a function of the other parameter. In Fig. 8a–e, for each parameter, the average of these percentages is presented for the test set.

The first parameter is the *area threshold*. In the circle-fit algorithm, the components smaller than this threshold are eliminated; thus, no circles smaller than this threshold exist in the image decomposition. Smaller values of this threshold result in representing noise as a set of objects. This decreases the segmentation performance. On the other hand, its larger values result in less number of circles, which causes to have some missing object information. This does not give a good representation of the image and decreases the performance too. To quantitatively understand the effect of this parameter, the experiments are repeated selecting the area threshold as {5, 10, 15, 20, 25, 30, 35, 40, 45, 50}. In Fig. 8a, the results are shown for the test set.

The second parameter is the *number of lumen neighbors* N . This parameter is used in initial seed determination, in which an object-graph is constructed defining edges between each lumen object and its N -closest lumen and N -closest nucleus objects. The local object-graph features of the lumen object are then used to classify it with either the gland or the non-gland class. The selection of smaller values of this parameter results in features lacking distinctive qualities, and hence, lowers segmentation performance. The selection of larger values does not increase the overall performance. However, when this parameter becomes very large, one has the risk of observing “curse of dimensionality”, which typically decreases the accuracy. With selecting this parameter as {1, 2, 3, 4, 5, 10, 15, 20, 25, 30}, the test results are shown in Fig. 8b.

The next parameter is the *number of nuclei neighbors* M , which is used to construct an object-graph to stop the region growing process. Smaller values of this parameter set less number of edges (barriers) at which region growing stops. This results in flooding of the regions and dramatically decreases the specificity, and thus, the segmentation accuracy. Larger values of this parameter do not change the accuracy too much. However, since more number of edges narrows down their inner regions, the final glands are smaller than expected for the larger values. This yields lower sensitivities but higher specificities. With selecting the number of nuclei neighbors as {1, 2, 3, 4, 5, 10, 15, 20, 25, 30}, the test results are shown in Fig. 8c.

The fourth parameter is the *small object threshold* P , which is used to eliminate smaller regions in gland seed growing. Here a region is eliminated if its area is smaller than the small object threshold percentage of the largest gland region in the image. Increasing this parameter results in eliminating more regions. This leads to higher specificity but lower sensitivity values. In the experiments, it is preferred to select a smaller value of this threshold to eliminate only the very small regions; false glands with larger sizes are eliminated in the false gland elimination step. The experiments are repeated selecting its value as {0.000, 0.010, 0.025, 0.050, 0.075, 0.100, 0.125, 0.150, 0.175, 0.200}. The results for the test set are shown in Fig. 8d.

The last parameter is the *simplification factor* k . It gives the degree of simplification of a polygon, which is formed by connecting the centroids of its nucleus objects. If this parameter is selected to be 1, there is no simplification for the polygon. If it is selected to be

⁷ The calculation of these rates is pixel-based. A pixel is considered as positive if it is identified as a gland pixel by the gland segmentation algorithm, and as negative otherwise. Therefore, the true positive rate (TP) is the number of positive pixels that belong to a gland in the gold standard; the false positive rate (FP) is the number of positive pixels that do not belong to a gland in the gold standard; the false negative rate (FN) is the number of negative pixels that belong to a gland in the gold standard; and the true negative rate (TN) is the number of negative pixels that do not belong to a gland in the gold standard.

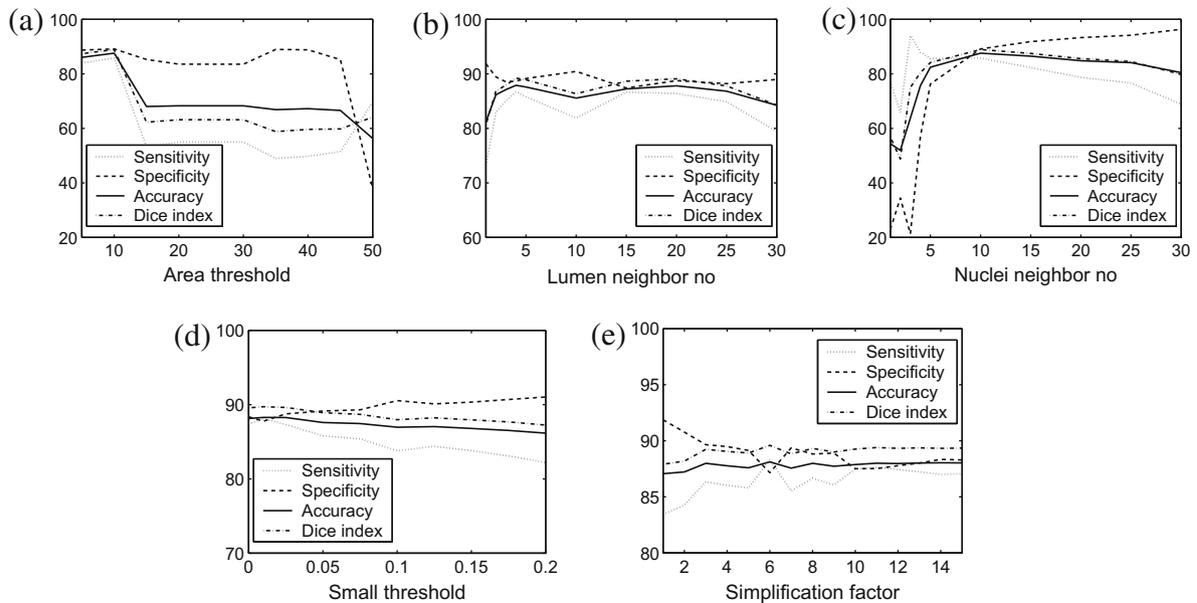


Fig. 8. For the test set, the sensitivity, specificity, accuracy, and Dice similarity index percentages as a function of (a) the area threshold, (b) the number of lumen neighbors N , (c) the number of nuclei neighbors M , (d) the small object threshold P , and (e) the simplification factor k .

the number of its nucleus objects, the simplification gives the convex hull of the polygon. This may not provide the correct gland borders since a gland is not necessarily convex. Thus, it should be selected in between these values. In Fig. 8e, the results are shown for the simplification factor being selected as {1,2,3,4,5,6,7,8,9,10,11,12,13,14,15} for the test set. These results show that the selection of this parameter affects the segmentation performance less compared to the other parameters.

3.2. Comparisons

This work compares the proposed algorithm, which relies on the use of object-graphs, with two previous approaches, both of which rely on the use of pixel-based information. In the first of these approaches (*nuclei-identification based approach*), the pixels of epithelial cell nuclei are first identified and the area surrounded by these identified pixels are then determined as in Wu et al. (2005a). In the second approach (*lumina-identification based approach*), the pixels of luminal areas are first identified and then the glands are obtained growing these identified pixels (Wu et al., 2005b). These algorithms are explained and their results are discussed in the following subsections.

3.2.1. The nuclei-identification based approach

The first step of the algorithm proposed in Wu et al. (2005a) is to distinguish the pixels of epithelial cell nuclei from the other pixels by thresholding the image. For this purpose, the algorithm first lowers the pixel intensities by applying a set of four directional filters (at the angles of 0° , 45° , 90° , and 135°) and selecting the lowest output of these filters for each pixel; each filter is defined as a 2-D Gaussian low-pass filter with standard deviations of σ_x and $\sigma_y = 4\sigma_x$. Then, the algorithm identifies the pixels with intensities smaller than an intensity threshold as the pixels of epithelial cell nuclei. Then, the algorithm dilates the identified nuclei pixels with a circular structuring element (with a radius of R) and fills the areas surrounded by these nuclei pixels. It finally identifies these areas as glands, provided that their areas are greater than an area threshold.

In Wu et al. (2005a) for the images that are taken with a $20\times$ microscope magnification and that have a resolution of 480×640 , the standard deviation σ_x and the radius R are selected

as 8 and 5, respectively. For each of these images, the intensity and area thresholds are manually selected. In our experiments, the value of these parameters is selected according to the segmentation performance obtained on the training set. In particular, all possible combinations of the following parameter sets are considered: a set of {16,32,...,224,240} plus the value obtained by the Otsu method (Otsu, 1979) for the intensity threshold; a set of {2,4,8,16} for the standard deviation σ_x ; a set of {3,5,7,9} for the radius R ; and a set of {0,500,1000,2500,5000,7500,10,000,12,500,15,000} pixels for the area threshold. Considering the parameter set that leads to the best sensitivity–specificity pair on the training samples, the intensity threshold is selected as 144, the standard deviation σ_x as 4, the radius R as 3, and the area threshold as 0.

3.2.2. The lumina-identification based approach

The second algorithm (Wu et al., 2005b) first thresholds the image to identify the nucleus and lumen pixels. Then, it determines the connected components of lumen pixels on which a round window (with a radius of R_0) can be located as initial gland seeds. Next, these initial gland seeds are iteratively grown dilating them with another round window (with a radius of R_i); a seed pixel is dilated if the round window centered at this pixel consists of only the lumen pixels. At the end of this iterative region growing process, the seeds for which the growth does not converge after a maximum number of iterations are considered as false glands and they are eliminated. Each remaining seed is dilated with a round window (with a radius of E) to find its surrounding dam and this seed is eliminated if the thickness of the dam (the ratio of its nucleus pixels) is smaller than a thickness threshold. Subsequently, nucleus pixels of the true glands are iteratively grown with a square structuring element (with a size of a) and final gland boundaries are obtained dilating these grown regions with a round window (with a radius of E_2).

In Wu et al. (2005b), for the images that are also taken with a $20\times$ microscope magnification and that have a resolution of 480×640 , the radii R_0 , R_i , E , and E_2 are selected as 25, 4, 10, and 10, respectively. The square size a is selected as 3 and the maximum number of iterations is selected as 50. Similar to the previous approach, the intensity and thickness thresholds are manually selected for each image. In our experiments, the value of these

parameters is selected according to the segmentation performance obtained on the training set. Similarly, all possible combinations of the following parameter sets are considered: a set of $\{16, 32, \dots, 224, 240\}$ plus the value obtained by the Otsu method for the intensity threshold; a set of $\{10, 25, 50\}$ for the radius R_0 ; a set of $\{2, 4, 8\}$ for the radius R_i ; a set of $\{5, 10, 15\}$ for the radius E ; a set of $\{5, 10, 15\}$ for the radius E_2 ; a set of $\{3, 5, 7\}$ for the square size a , a set of $\{2.0, 4.0, 6.0, 8.0, 10.0\}$ for the thickness threshold. The maximum number of iterations is also selected as 50. Considering the parameter set that leads to the best sensitivity–specificity pair on the training samples, the intensity threshold is selected as 112, the radius R_0 as 25, the radius R_i as 8, the radius E as 10, the radius E_2 as 5, the square size a as 7, and the thickness threshold as 2.0.

3.2.3. Results

For the training and test sets, the sensitivity, specificity, accuracy, and Dice similarity index of the nuclei-identification and lumina-identification based approaches with the selected

parameter sets are reported in Tables 3 and 4. These tables also present the object-graph results that are obtained both before and after applying the false gland elimination step (FGE step). The results demonstrate that the object-graph approach, both with and without false gland elimination, improves the segmentation accuracy as well as the Dice similarity index of the pixel-based segmentation approaches, leading to high sensitivity and specificity values. The Wilcoxon test with a significance level of 0.05 exhibits that this improvement is statistically significant.

For both the nuclei-identification and lumina-identification based approaches, the visual segmentation results for the images given in Fig. 2 are also illustrated in Figs. 9 and 10, respectively. These figures show that good segmentations are obtained for only a few images. They also show that some segmentation results are inconsistent with our previous intuitions, which are discussed in the introduction; for example, it is expected that the nuclei-identification based approach yields better segmentation results for the image given in Fig. 2g since the nucleus pixels of this image are expected to form closed components. Thus, for this image, the other

Table 3
For the object-graph approach, the nuclei-identification based approach, and the lumina-identification based approach, the average and the standard deviation of the sensitivity, specificity, accuracy, and Dice similarity index percentages obtained on the *training* set.

	Sensitivity	Specificity	Accuracy	Dice index
Object-graphs (before FGE step)	89.61 ± 4.28	65.09 ± 25.44	77.39 ± 14.99	81.73 ± 12.71
Object-graphs (after FGE step)	83.43 ± 7.73	92.30 ± 5.78	88.00 ± 4.16	88.46 ± 4.62
Nuclei-identification	55.88 ± 28.48	55.16 ± 32.47	56.34 ± 18.31	57.27 ± 19.71
Lumina-identification	47.24 ± 29.81	92.70 ± 8.42	68.58 ± 12.75	56.12 ± 30.54

Table 4
For the object-graph approach, the nuclei-identification based approach, and the lumina-identification based approach, the average and the standard deviation of the sensitivity, specificity, accuracy, and Dice similarity index percentages obtained on the *test* set.

	Sensitivity	Specificity	Accuracy	Dice index
Object-graphs (before FGE step)	90.62 ± 5.44	72.80 ± 15.38	82.57 ± 8.36	84.31 ± 9.76
Object-graphs (after FGE step)	85.80 ± 6.71	89.14 ± 10.40	87.59 ± 5.01	88.91 ± 4.63
Nuclei-identification	53.77 ± 25.67	51.67 ± 33.64	53.24 ± 13.62	54.33 ± 19.69
Lumina-identification	52.59 ± 32.88	87.48 ± 15.12	67.62 ± 17.17	59.04 ± 30.00

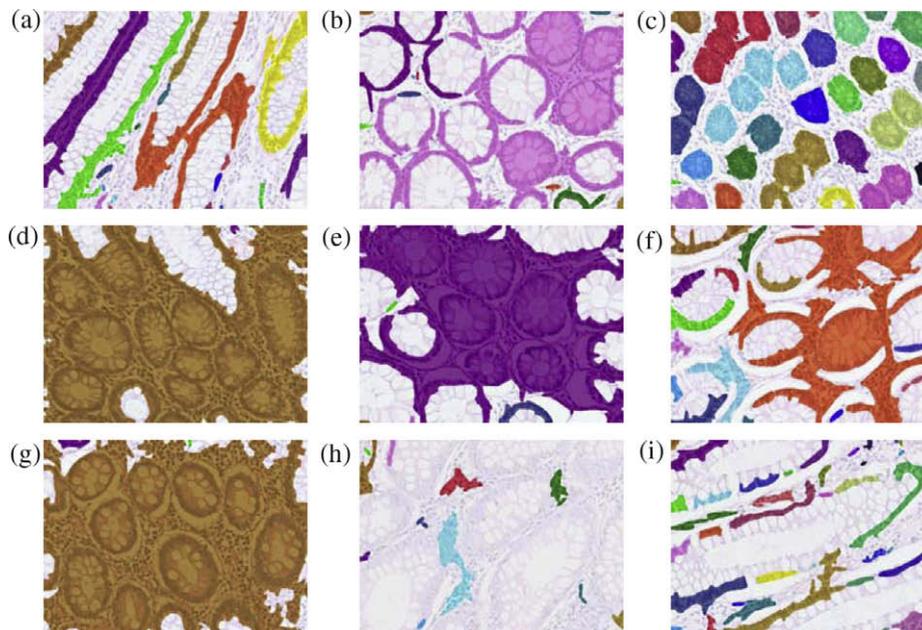


Fig. 9. The visual results obtained by the nuclei-identification based approach for the tissue images given in Fig. 2.

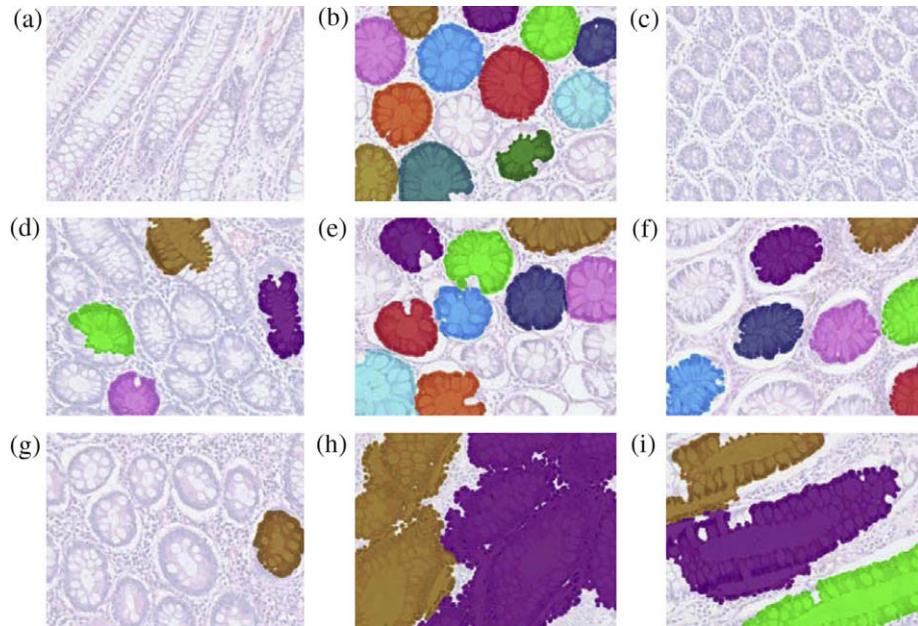


Fig. 10. The visual results obtained by the lumina-identification based approach for the tissue images given in Fig. 2.

segmentation results in which different parameter sets are used are also examined. Here it is observed that better segmentation results could be obtained when the parameter set is optimized according to this image. However, this decreases the segmentation performance for the others. This is also true for the other images. When the parameter set is optimized for a particular image, the segmentation performance of the two pixel-based approaches (especially those of the lumina-identification based approach) could increase for this particular image. However, the segmentation performances decrease for the others. This demonstrates the difficulty of selecting a single parameter set (for the pixel-based approaches) that would work for all images.

4. Discussions

In this work, we introduce an object-based approach for the purpose of gland segmentation. This approach decomposes the tissue image into a set of primitive objects and segments glands making use of the spatial distributions of these objects, which are quantified with the definition of object-graphs. In this work, the experiments are conducted on the images of 72 colon tissues of 36 different patients. Experimental results demonstrate that the proposed object-based approach yields high accuracies of 77.39% for the training set and 82.57% for the test set and significantly improves the segmentation performance of its pixel-based counterparts. Furthermore, with a false gland elimination step, these accuracies increase up to 88.00% for the training set and 87.59%

for the test set. These results show that the use of object-based information, instead of using pixel-based information alone, leads to more robust segmentations to imaging artifacts. This is attributed to pixel intensities being more sensitive to the noise that arises from the staining, fixation, and sectioning related problems.

The proposed method is expected not to be sensitive with respect to a particular microscope and to also work with images that are taken from other microscopes. To examine this issue, some images are taken with a camera mounted onto an Olympus BX51 Microscope. On these images, the proposed algorithm is run without modifying any of its parameters and without changing the rules of its decision tree. The segmentation results obtained on these samples show that the proposed method could also work for other microscopes. These segmentation results are visually illustrated in Fig. 11.

The proposed algorithm provides an infrastructure for further analysis of biopsies that include glandular structures. This infrastructure allows us to locate glands on the tissue image and to understand whether or not a gland deviates from its normal structure. To identify such glands, one method would be to extract a set of mathematical features from each of the segmented glands and to classify the glands using these mathematical features. Another method would be to quantify the spatial relations of the segmented glands extracting structural features; e.g., graphs or Voronoi diagrams could be defined considering each gland as a node and then local features extracted for each node could be used to classify the gland. Similarly, global features extracted for the entire graph or the Voronoi diagram could be used to classify the entire

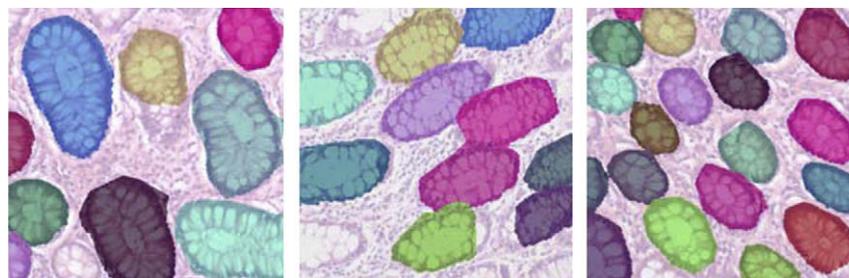


Fig. 11. The visual results for the histopathological images that are taken with a camera mounted onto an Olympus BX51 Microscope.

tissue. As a future research work, we plan to investigate such features for the purpose of gland classification.

Acknowledgment

This work has been supported by the Scientific and Technological Research Council of Turkey under the project number TÜBİTAK 106E118.

References

- Ahmadyfard, A., Kittler, J., 2003. A multiple classifier system approach to affine invariant object recognition. In: Proceedings of the International Conference on Computer Vision Systems, Graz, Austria.
- Andrion, A., Magnani, C., Betta, P.G., Donna, A., Mollo, F., Scelsi, M., Bernardi, P., Botta, M., Terracini, B., 1995. Malignant mesothelioma of the pleura: interobserver variability. *Journal of Clinical Pathology* 48, 856–860.
- Choi, H.-K., Jarkrans, T., Bengtsson, E., Vasko, J., Wester, K., Malmstrom, P.-U., Busch, C., 1997. Image analysis based grading of bladder carcinoma. Comparison of object, texture and graph based methods and their reproducibility. *Analytical Cellular Pathology* 15, 1–18.
- Christmas, W.J., Kittler, J., Petrou, M., 1995. Structural matching in computer vision using probabilistic relaxation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 17 (8), 749–764.
- Demir, C., Gultekin, S.H., Yener, B., 2005. Learning the topological properties of brain tumors. *IEEE-ACM Transactions on Computational Biology and Bioinformatics* 2 (3), 262–270.
- Esgiar, A.N., Naguib, R.N.G., Sharif, B.S., Bennett, M.K., Murray, A., 1998. Microscopic image analysis for quantitative measurement and feature identification of normal and cancerous colonic mucosa. *IEEE Transactions on Information Technology in Biomedicine* 6, 197–203.
- Esgiar, A.N., Naguib, R.N.G., Sharif, B.S., Bennett, M.K., Murray, A., 2002. Fractal analysis in the detection of colonic cancer images. *IEEE Transactions on Information Technology in Biomedicine* 6, 54–58.
- Farjam, R., Soltanian-Zadeh, H., Jafari-Khouzani, K., Zoroofi, R.A., 2007. An image analysis approach for automatic malignancy determination of prostate pathological images. *Clinical Cytometry* 72B (4), 227–240.
- Gunduz-Demir, C., 2007. Mathematical modeling of the malignancy of cancer using graph evolution. *Mathematical Biosciences* 209 (2), 514–527.
- Hamilton, P.W., Bartels, P.H., Thompson, D., Anderson, N.H., Montironi, R., 1997. Automated location of dysplastic fields in colorectal histology using image texture analysis. *Journal of Pathology* 182, 68–75.
- Jain, B.J., Wysotzki, F., 2004. Central clustering of attributed graphs. *Machine Learning* 56 (1–3), 169–207.
- Keenan, S.J., Diamond, J., McCluggage, W.G., Bharucha, H., Thompson, D., Bartels, B.H., Hamilton, P.W., 2000. An automated machine vision system for the histological grading of cervical intraepithelial neoplasia (CIN). *Journal of Pathology* 192, 351–362.
- Luo, B., Wilson, R.C., Hancock, E.R., 2006. A spectral approach to learning structural variations in graphs. *Pattern Recognition* 39 (6), 1188–1198.
- Naik, S., Doyle, S., Feldman, M., Tomaszewski, J., Madabhushi, A., 2007. Gland segmentation and Gleason grading of prostate histology by integrating low-, high-level and domain specific information. In: Proceedings of the 2nd Workshop on Microscopic Image Analysis with Applications in Biology, Piscataway, NJ.
- Nielsen, B., Albrechtsen, F., Danielsen, H.E., 1999. The use of fractal features from the periphery of cell nuclei as a classification tool. *Analytical Cellular Pathology* 19, 21–37.
- Otsu, N., 1979. A threshold selection method from gray level histograms. *IEEE Transactions on Systems Man and Cybernetics* 9, 62–66.
- Petrakis, E.G.M., Faloutsos, C., 1997. Similarity searching in medical image databases. *IEEE Transactions on Knowledge and Data Engineering* 9 (3), 435–447.
- Sanfeliu, A., Alquezar, R., Andrade, J., Climent, J., Serratos, F., Verges, J., 2002. Graph-based representations and techniques for image processing and image analysis. *Pattern Recognition* 35 (3), 639–650.
- Serra, J., 1982. *Image Analysis and Mathematical Morphology*. Academic, London.
- Spyridonos, P., Ravazoula, P., Cavouras, D., Berberidis, K., Nikiforidis, G., 2001. Computer-based grading of haematoxylin–eosin stained tissue sections of urinary bladder carcinomas. *Medical Informatics and the Internet in Medicine* 26, 179–190.
- Tang, F., Tao, H., 2008. Probabilistic object tracking with dynamic attributed relational feature graph. *IEEE Transactions on Circuits and Systems for Video Technology* 18 (8), 1064–1074.
- Thiran, J.-P., Macq, B., 1996. Morphological feature extraction for the classification of digital images of cancerous tissues. *IEEE Transactions on Biomedical Engineering* 43, 1011–1020.
- Thomas, G.D., Dixon, M.F., Smeeton, N.C., Williams, N.S., 1983. Observer variation in the histological grading of rectal carcinoma. *Journal of Clinical Pathology* 36, 385–391.
- Tosun, A.B., Kandemir, M., Sokmensuer, C., Gunduz-Demir, C., 2009. Object-oriented texture analysis for the unsupervised segmentation of biopsy images for cancer detection. *Pattern Recognition* 42 (6), 1104–1112.
- Weyn, B., Van de Wouwer, G., Kumar-Singh, S., Van Daele, A., Scheunders, P., Van Marck, E., Jacob, W., 1999. Computer-assisted differential diagnosis of malignant mesothelioma based on syntactic structure analysis. *Cytometry* 35, 23–29.
- Wiltgen, M., Gerger, A., Smolle, J., 2003. Tissue counter analysis of benign common nevi and malignant melanoma. *International Journal of Medical Informatics* 69, 17–28.
- Wolberg, W.H., Street, W.N., Heisey, D.M., Mangasarian, O.L., 1995. Computer-derived nuclear features distinguish malignant from benign breast cytology. *Human Pathology* 26, 792–796.
- Wu, H.-S., Xu, R., Harpaz, N., Burstein, D., Gil, J., 2005a. Segmentation of microscopic images of small intestinal glands with directional 2-D filters. *Analytical and Quantitative Cytology and Histology* 27 (5), 291–300.
- Wu, H.-S., Xu, R., Harpaz, N., Burstein, D., Gil, J., 2005b. Segmentation of intestinal gland images with iterative region growing. *Journal of Microscopy* 220 (3), 190–204.