

# AN EFFICIENT METHOD FOR DETECTION AND SEGMENTATION OF NUCLEI IN DIGITAL PATHOLOGY

VENAKTA KARTEEK PALADUGU AND THOMAS B. KINSMAN, PH.D.

## 1. OVERVIEW

Pathology is the core in any medical procedure[6]. Traditionally humans analyze numerous biopsies in order to recognize any anomalies in the tissue and this is a very intensive task, but Digital Pathology has been very helpful in minimizing human intervention and helping the pathologist to minimize the work and assist him. Histopathology is examining a piece of tissue under a microscope through thinly slicing and staining the tissue sample with dyes. It provides a very accurate information about tissues and can be useful in classifying cells to be benign or malignant. This paper provides an effective method for detecting nuclei and segment them in digital histopathology imagery[16].

Detecting Nuclei and segmenting them are important in cancer diagnosis. The following are the most commonly used methods in Imaging chain for detecting nuclei and segmenting them:

- (1) *Thresholding*
- (2) *Morphology*
- (3) *Region Growing*
- (4) *Watershed*
- (5) *Active Contour Models and Level sets*
- (6) *K-means clustering*
- (7) *Probabilistic models*
- (8) *Graph cuts*

First preprocessing such as thresholding[11], morphology[14] and histogram equalization is done on the pathology images for noise reduction. Secondly nuclei are detected using multiscale Laplacian of Gaussian (LoG) filtering along with Euclidean distance map of the binarized image .Thirdly nuclei segmentation which is based on size, texture, shape and other morphological appearances is done using thresholding and morphological operations. Finally, classification is done using features precomputed from image segmentation which is totally based on the quality of segmentation.

Many of the studies have been done using their own dataset so in order to compare different studies there should be a benchmark dataset. These datasets should be validated by different pathologists and must be taken from many patients worldwide. Many of the methods above can detect the region of the nuclei but a lot of difficulties arise during detection of overlapping nuclei, so this field is still an open research area. It is also important to address the robustness of different clinical/technical conditions such as different apparatus used for obtaining the image, various characters of stains, various conditions in which the image is obtained.

## 2. IMAGE PROCESSING METHODS

A.Thresholding: Thresholding can be used to remove the outliers from the image. We can find the threshold value using Otsu method which determines the threshold that minimizes the intra class variance [9]. There are other methods like adaptive thresholding which finds thresholds based on regions so thresholds may vary across the image this can be used to remove non uniform illumination.

**B.Morphology:** There are 4 types of morphological operations that can be used for noise removal or for edge detection. The basic operations are dilution and erosion. Dilation can be used to enlarge the foreground pixels while erosion can be used to shrink them. The other two are combination of the above two basic operations namely Closing and Opening. Closing is a combined operation of the two basic operations dilation of an image followed by erosion doing so will remove holes in the object of the image, whereas Opening is also similar to Closing but it does erosion first and then it does dilation which removes very small dots in the image. These 4 operations are used to remove salt and pepper noise.

**C.Region Growing:** In this method we first choose the seed points and then grow regions from them based on their similarity to the seed points[17]. Selection of seed points can be done by computing the histogram and choosing the threshold value having good amount of pixels i.e. high peaks. The homogeneity predicates often chosen are average intensity, variance color, texture, motion, shape and size.

**D.Watershed:** This can be used if two regions of interest are very close to each other that is their edges are touching one another. This technique treats the image as a topographic map where the heights are the intensity levels of each pixel[13]. For example, low intensity values can be of lower height and high intensity values can be hills or mountain ridges. Catchment basins are the areas with lowest intensity, and we put a source of water at each basin and start flooding then the areas where the floodwater from different basins meet are identified and a barrier is formed within the pixels in these areas.

**E.Active Contour Models and Level sets:** Active contour models can be used to obtain deformable models in image segmentation[5]. Curvature of the model is determined by minimizing external and internal energies. External energy deals with the positioning of the contour in the image and internal energy deals with controlling the deformable changes like elasticity and stiffness of the contour. These energies are defined by the points on the curve and their derivative along the curve, when these energies are minimum that means the contour has converged resulting in the object contour that we required.

**F.K-means Clustering:** It is an iterative process, the algorithm is as follows[7]:

- (1) Randomly selecting k cluster centers.
- (2) Assign each pixel to a cluster based on minimum distance among all centers.
- (3) Average all the pixels in each cluster and reassign the center to be the average.
- (4) Keep repeating 2nd and 3rd step until convergence is attained i.e. the centers do not change after certain point.

These methods are generally used in the Imaging chain of detecting nuclei.

### 3. ARCHITECTURE

The architecture is as follows:

- (1) Firstly the input image is taken and preprocessing is done to remove the noise.
- (2) Secondly nuclei detection is done and then based on that we do the ellipse detector analysis.
- (3) Ellipse detector finds whether a nucleus is isolated or overlapping with another so that we can do segmentation.
- (4) Finally we segment accordingly based on whether the nuclei is isolated or not.

### 4. PREPROCESSING

First step before nuclei segmentation is removing the noise from the image and there are various methods that can be used as mentioned below:

- (1) **Illumination Normalization:** This is used to remove the unevenness in the illumination. The unevenness in the illumination can be normalized either by referring through all images and getting the average luminance and changing each image according to the average or through white shading

correction. In white shading correction we correct the illumination by taking a blank image and it is used to correct the original image pixel by pixel by the following formulae[8]:

$$(1) \quad Transmittance = \frac{SpecimenValue - BackgroundValue}{WhiteReferenceValue - BackgroundValue}$$

- (2) Noise Reduction: Using thresholding, pixels that lie outside the threshold are removed. The threshold can be found out by observing histogram and the outliers threshold value. This removes unwanted noise and then we can use morphological operations like opening and closing to remove noise while preserving the nuclei shape[4]. We can also use histogram equalization to improve the contrast between foreground and background regions and gaussian filtering can be used to smoothen out the nuclei regions.

## 5. NUCLEI DETECTION

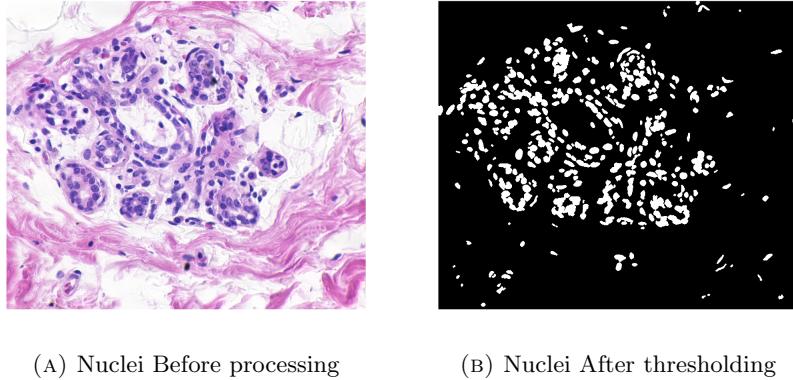


FIGURE 1. Nuclei before and after thresholding

For the nuclei detection we are going to use adaptive thresholding and the procedure is as follows[16]:

- (1) Firstly to segment the background from foreground looked through all channels and the best contrast was given by Red channel of RGB. Later opening by reconstruction and closing by reconstruction are done to remove noise and make the regions of nuclei homogeneous.
- (2) Secondly combination of local and global thresholding is applied on the image. After finding a threshold we convert the above image into binary image by setting values less than threshold to 0 and vice versa.
- (3) Thirdly the binarized image is further processed to remove small connected components and also to make regions of nuclei homogeneous using morphology.

The resulting image after the above steps is shown in figure 1(B). In the later steps we use the Ellipse Descriptor Analysis to check whether a nucleus is isolated or not and proceed accordingly.

## 6. ELLIPSE DESCRIPTOR ANALYSIS

After the above segmentation is done the nuclei can be isolated or clumped. In this section isolated nuclei are separated from clumped nuclei using ellipses, based on the observations of Hongming Xu et al.[16] the shape of single nuclei is generally elliptical. In order to build this model, we need to get the boundary points of all connected components and that can be done through getting the points with non-zero gradient magnitudes. As the points within the nuclei do not differ along x or y and only the border points have sudden changes, these points will have non zero gradient magnitudes.

After finding boundary points an ellipse is tried to fit along the boundary points using the fitellipse[3] function in MATLAB. After fitting the ellipse, evaluation of the points within the ellipse is done to find whether the nuclei is clumped or isolated, this can be done through as follows, let  $P_n$  be the list of pixels

present in nuclei and  $P_e$  be the list of pixels that are within the ellipse. The following two parameters are defined to measure the ellipticity of the region:

$$(2) \quad e_1 = \frac{|P_n \cap P_e|}{|P_e|}$$

$$(3) \quad e_2 = \frac{|P_n \Delta P_e|}{|P_e|}$$

If  $e_1$  is higher i.e. close to 1 and  $e_2$  is lower i.e. close to 0 then that means ellipse fit correctly on the nuclei which means it is isolated whereas if  $e_1$  is lower and  $e_2$  is higher then that means ellipse didn't fit the nuclei correctly and thus it is clumped. Based on the research done on the nuclei regions in the nuclei segmentation done by Hongming Xu et al.[16] the threshold values for  $e_1$  and  $e_2$  for isolated nuclei are considered as 0.91 and 0.18. So if  $e_1 > 0.91$  and  $e_2 < 0.18$  then it is considered as isolated nuclei and segmented accordingly. In the next section a voting algorithm is described in order to segment clumped nuclei.

## 7. CONE VOTING ALGORITHM

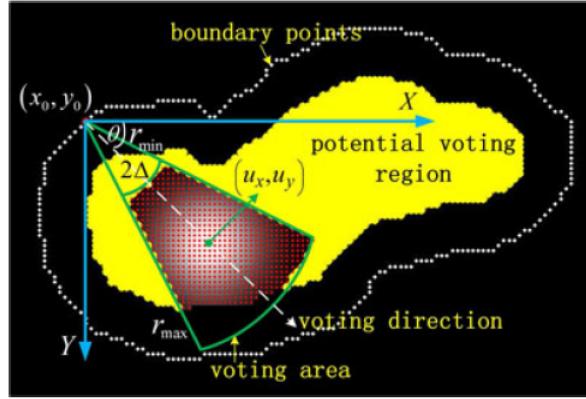


FIGURE 2. Illustration of modified cone voting area[16].

In order to segment the clumped nuclei we need to determine the seed points i.e. nuclei centers within the clumped nuclei and later apply marked watershed algorithm using those seed points to segment. In this section using voting algorithm seed points are determined. Based on the Parvin et al. [10] MPV algorithm and Qi et al. [12] SPV algorithm this voting algorithm with few modifications is designed. A cone is formed from each boundary point of the clumped nuclei with two radii  $r_{min}$  and  $r_{max}$  based on the observations of nuclei sizes and angle  $\Delta$ [16] as shown in figure 2. Pixels inside the cone get voted based on how far they are from the Gaussian kernel centred at core of the cone i.e points close to centre of the core get more votes. The voting also depends on the gradient magnitude of the border point.

The Guassian kernel is defined so that the center point in the cone gets more weightage and so the seed point which is near the center gets more votes. Before forming a cone, erosion on all clumped nuclei is done so that points close to the border will not be considered while voting as the nuclei center will not be present there. The computational time is improved because of the erosion operation as there are less points while voting and the Gaussian kernel ensures false seeds detection. The voting equation for each point is as follows:

$$(4) \quad V(x, y) = V(x, y) + \sum_{(x,y) \in A(m,n)} ||\nabla T(m, n)|| g(x, y, u_m, u_n, \sigma)$$

where  $V(x,y)$  is the voting image of the same resolution as original image,  $(x,y)$  represents coordinates of the pixels and  $A(m,n)$  represents the voting area for the border point  $(m,n)$ ,  $\nabla T(m,n)$  is the gradient magnitude of the border point  $(m,n)$  and  $g(x,y,u_m,u_n, \sigma)$  is the Gaussian kernel centred at  $(u_m, u_n)$  defined as below:

$$(5) \quad g(x, y, u_m, u_n, \sigma) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{(x - u_m)^2 + (y - u_n)^2}{2\sigma^2}\right)$$

where  $u_m = m + r\cos(\theta)$ ,  $u_n = n + r\sin(\theta)$  and  $r = \frac{(r_{min}+r_{max})}{2}$ . In the next sections detecting seed points from voting image and marked watershed algorithm are described.

## 8. SEED DETECTION

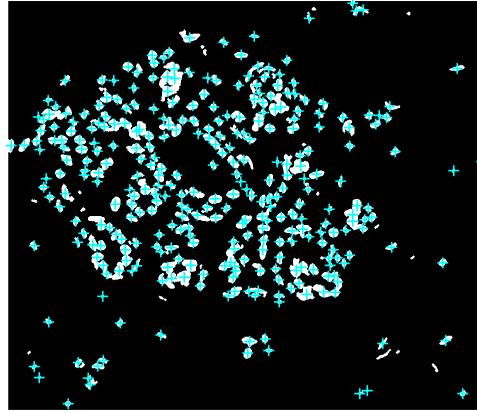


FIGURE 3. Seed points detected from voting image after Mean shift algorithm

From the voting image we will do mean shift clustering which will give us the mode pixel of each component which are our seed points. The mean shift clustering is done as follows [2]:

- (1) A point is considered randomly from all the voting pixels that are calculated in previous section.
- (2) Then all points within its neighborhood i.e. based on average nuclei size within a bandwidth of 7.2 units in terms of Euclidean distance are considered.
- (3) After getting the neighborhood points, mean of all the neighborhood points is calculated with a linear kernel but a gaussian kernel can also be used and shift the old point to the new mean.
- (4) 3<sup>rd</sup> step is iterated continuously until the mean converges i.e. mean does not shift from its previous value.
- (5) The above converged mean value is considered as the seed point for one component.
- (6) 1<sup>st</sup> step is repeated by not considering the points visited in 2<sup>nd</sup> to 4<sup>th</sup> steps and it continues until there are no more points to visit.

The resulting seed points are shown in Figure 3. After doing the mean shift clustering, the seed points are collected and then marked watershed algorithm is applied which is described in next section.

## 9. NUCLEI SEGMENTATION

For segmentation of clumped nuclei general watershed segmentation can be used which considers inner distance transform of the image as a topological surface and finds minima in each region and floods the surface around the minima called catchment basin[15]. After flooding water merging from two minima are prevented by forming a line called watershed lines which are nothing but lines that separate catchment basins and in the nuclei segmentation perspective, they are the lines that separate the overlapping nuclei. But this results in over segmentation as there can be many minima within a region and that's where seed points are useful.

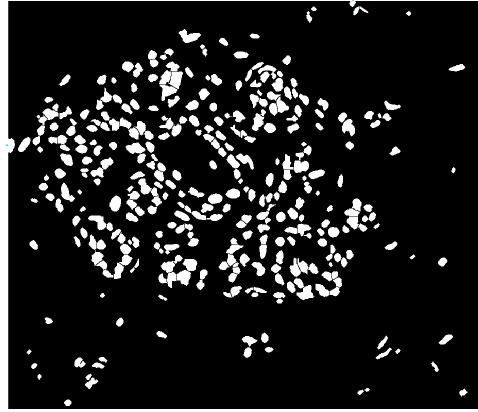


FIGURE 4. Segmentation of clumped nuclei after Marked Watershed Algorithm

Using seed points as minima we do watershed segmentation which is called as marked watershed segmentation. But the classic marker function [1] which is based on inner distance transform fails to separate highly overlapping nuclei causing under segmentation, so a new marker function which takes account of inner distance transform and the Euclidean distance is calculated[16].The new marker function for the image containing nuclei regions and  $C = (C_1, C_2, \dots, C_k)$  be the detected nuclei seeds is as follows:

$$(6) \quad f(x, y) = \min_j \{D_j(x, y) - d_i(C_j)\}$$

Where  $D_j(x, y)$  is the Euclidean distance between jth seed  $C_j(1 \leq j \leq k)$  and the pixel  $(x,y)$  and  $d$  is the inner distance map for the pixel . This new marker function will stop under segmentation because of consideration of two distances and help in forming a better water line. The resulting image after marked watershed algorithm using new marker function is show in Figure 4.

## 10. RESULTS

The following Figure 5 shows the initial and final results of 5 images of nuclei.

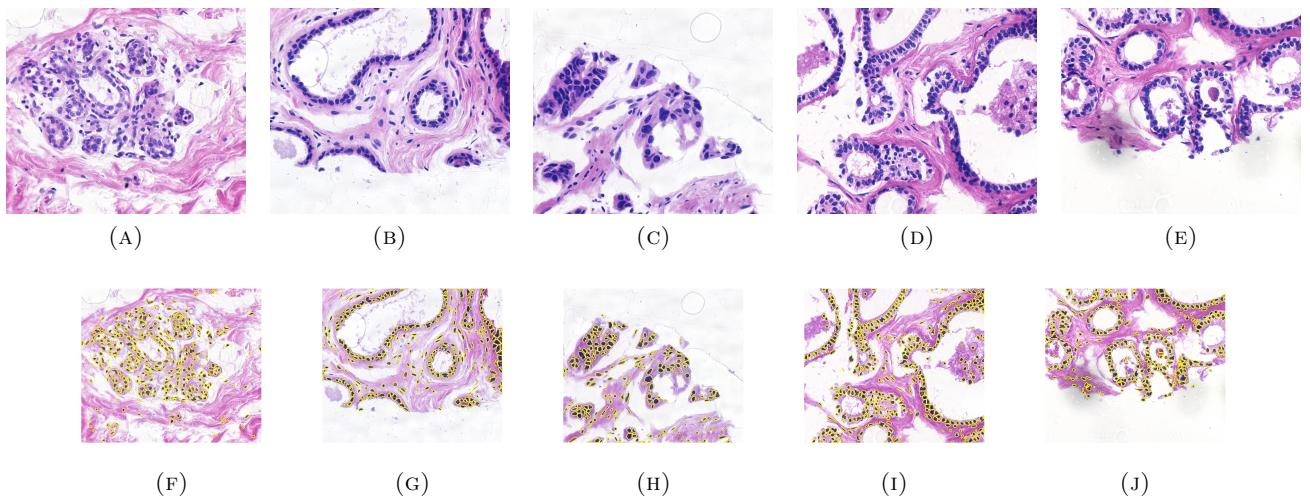


FIGURE 5. Figures (A)-(E) represent nuclei before detection and Figures (F)-(J) represent nuclei after segmentation and detection.

## 11. LESSONS LEARNED

Firstly pre-processing is very important as noise can decrease accuracy, so morphological operations can be used to remove the noise. Shapes can tell us a lot about image,taking the advantage of it separated isolated and clumped nuclei based on ellipse shape.Secondly learnt a voting algorithm that can be used to detect seed points in the clumped regions and the importance of weights in the voting that help in detecting the seed point accurately, in this case it was Guassian centred at center of cone but we can alter it depending on where we want more weight to be at.

Thirdly learnt Mean shift clustering which works by selecting one or more random points and gives points where most of the data is around i.e. mode using mean of the data.Finally learnt about two important segmentation called Watershed and Marked Watershed Algorithm which segment the nuclei by filling the catchment basins from the minima and form water lines when two flows from different basins meet, the difference between them both is Watershed considers all minima in the image but the Marked Watershed Algorithm considers only the given seed points as minima and forms catchment basins around them. Also learnt the use of distance transforms that help in forming the toplogical surface required for the Watershed algorithms.

## 12. FUTURE DIRECTIONS

One of the challenges in this field is lack of unified benchmarks. Research performed in this area are based on local datasets that are with the researches.This is bad because different studies have different datasets and since they are not tested on a large scale many inaccuracies may come. So benchmark datasets compromising of large number of patients validated by different pathologists should be set.Other challenge includes ellipse fitting techniques which sometimes cannot accurately fit shape of nuclei so detection of nuclei is not yet perfect therefore more research should be done here. Finally the robustness of various clinical conditions should be addressed such as apparatus used for obtaining the image, different lightening conditions and different staining characteristics as it may cause irregularites in the image resulting in inaccuracy.

## 13. CONCLUSION

This paper provides efficient methods for nuclei detection and segmentation. Initially the Red color channel is chosen as it looked better for segmentation among all channels. Later adaptive thresholding is applied in order to segment background from foreground and then morphological operations are applied in order to remove small noise leftover afterwards. After segmentation isolated nuclei can be detected easily as they fit the shape of ellipse but clumped nuclei are difficult to detect.

For the detection of clumped nuclei a voting algorithm based on the Parvin et al. [10] MPV algorithm and Qi et al. [12] SPV algorithm is implemented and mean shift clustering is done on the voting image to get the nuclei seeds for each componenet. After finding the seed points marked watershed algorithm with a custom marker function is used to do segmentation of clumped nuclei.

Finally because of the two improvements in the voting algorithm i.e erosion for lesser points while voting and Gaussian weight to ensure detection of false seed points improved the time complexity and accuracy of nuclei detection.

## 14. ACKNOWLEDGEMENT

I am grateful to Dr. Thomas B. Kinsman for giving me the opportunity to explore an area of my interest and for reviewing and guiding me in the project throughout.

## REFERENCES

- [1] J. Cheng and J. C. Rajapakse \*. Segmentation of clustered nuclei with shape markers and marking function. *IEEE Transactions on Biomedical Engineering*, 56(3):741–748, 2009.

- [2] D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Trans. Pattern Anal. Mach. Intell.*, 24:603–619, 2002.
- [3] A. Fitzgibbon, M. Pilu, and R. B. Fisher. Direct least square fitting of ellipses. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(5):476–480, 1999.
- [4] P. W. HAMILTON, P. H. BARTELS, D. THOMPSON, N. H. ANDERSON, R. MONTIRONI, and J. M. SLOAN. Automated location of dysplastic fields in colorectal histology using image texture analysis. *The Journal of Pathology*, 182(1):68–75, 1997.
- [5] M. Kass, A. Witkin, and D. Terzopoulos. Snakes: Active contour models. *INTERNATIONAL JOURNAL OF COMPUTER VISION*, 1(4):321–331, 1988.
- [6] W. W. Ma and A. A. Adjei. Novel agents on the horizon for cancer therapy. *CA: A Cancer Journal for Clinicians*, 59(2):111–137, 2009.
- [7] J. Macqueen. Some methods for classification and analysis of multivariate observations. pages 281–297, 1967.
- [8] G. Marty. Blank-field correction for achieving a uniform white background in brightfield digital photomicrographs. *BioTechniques*, 42:716, 718, 720, 07 2007.
- [9] N. Otsu. *An automatic threshold selection method based on discriminant and least squares criteria*. 1979.
- [10] B. Parvin, Q. Yang, J. Han, H. Chang, B. Rydberg, and M. H. Barcellos-Hoff. Iterative voting for inference of structural saliency and characterization of subcellular events. *IEEE Transactions on Image Processing*, 16(3):615–623, 2007.
- [11] A. Păsărică, R. G. Bozomitu, O. Diana Eva, D. Tărniceriu, and C. Rotariu. Analysis of different threshold selection methods for eye image segmentation used in eye tracking applications. pages 299–302, 2016.
- [12] X. Qi, F. Xing, D. J. Foran, and L. Yang. Robust segmentation of overlapping cells in histopathology specimens using parallel seed detection and repulsive level set. *IEEE Transactions on Biomedical Engineering*, 59(3):754–765, 2012.
- [13] J. Roerdink and A. Meijster. The watershed transform: Definitions, algorithms and parallelization strategies. *Fundam Inf*, 41, 10 2003.
- [14] J. Serra. *Image Analysis and Mathematical Morphology*. Academic Press, Inc., USA, 1983.
- [15] L. Vincent and P. Soille. Watersheds in digital spaces: an efficient algorithm based on immersion simulations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(6):583–598, 1991.
- [16] H. Xu, C. Lu, and M. K. Mandal. An efficient technique for nuclei segmentation based on ellipse descriptor analysis and improved seed detection algorithm. *IEEE Journal of Biomedical and Health Informatics*, 18:1729–1741, 2014.
- [17] S. W. Zucker. Region growing: Childhood and adolescence\*. *Computer Graphics and Image Processing*, 5:382–399, 1976.