

# **Improving Cancer Subtype Diagnosis and Grading using Clinical Decision Support System Based on Computer-Aided Tissue Image Analysis**

A Dissertation  
Presented to  
The Academic Faculty

by

Qaiser Mahmood Chaudry

In Partial Fulfillment  
of the Requirements for the Degree  
PhD in the  
School of Electrical and Computer Engineering

Georgia Institute of Technology  
May 2013

# **Improving Cancer Subtype Diagnosis and Grading using Clinical Decision Support System Based on Computer-Aided Tissue Image Analysis**

Approved by:

Dr. May D. Wang, Advisor  
Wallace H. Coulter Dept. of Biomedical  
Engineering  
*Georgia Institute of Technology*

Dr. Edward J. Coyle  
School of Electrical and Computer  
Engineering  
*Georgia Institute of Technology*

Dr. Paul Benkeser  
Wallace H. Coulter Dept. of Biomedical  
Engineering  
*Georgia Institute of Technology*

Dr. Andrew N. Young  
Pathology and Laboratory Medicine  
*Emory University School of Medicine,  
Atlanta, GA*

Dr. Anthony J. Yezzi  
School of Electrical and Computer  
Engineering  
*Georgia Institute of Technology*

Date Approved: Nov 19, 2012

To my parents, my wife Farzana, my son Aitazaz and my daughters Sundas and Iman

## ACKNOWLEDGEMENTS

*My deepest thanks to Allah for giving me courage and devotion to take all the challenges in my life; I have always been blessed by His mercy.*

I would like to express my gratitude to my thesis advisor, Professor May Dongmei Wang, for her guidance, support, and, help. She has been and will remain a source of motivation for me.

My profound thanks to the members of my reading committee Professors Paul Benkeser and Anthony Yezzi for their extremely valuable comments for shaping this thesis better than what it was. I would like to thank Professors Edward Coyle and Andrew Young for serving in my defense committee. I am especially grateful for the financial support provided to me by the Government of Pakistan and School of Electrical and Computer Engineering for this research.

My special thanks go to the Miblab team especially S. Hussain Raza, Yachna Sharma, Sonal Kothari, Dr. Richard Moffit and Dr. Jhon Phan, Dr. Mitchell Perry for their invaluable comments and suggestions for my research, thesis and presentation.

I would like to express my deepest gratitude to my wife, Farzana. She is an excellent wife and wonderful mother. Her love and faith in me has been a constant source of happiness in my life. My special love goes to my three children, Aitazaz, Sundas and Iman who think schooling never ends in your life as did I when I used to go together with my father to school(he being Principal of the school).

I am deeply grateful to my family and friends, especially to my parents for their love and encouragement throughout my education. All my accomplishments are due to their love and prayers.

# TABLE OF CONTENTS

CHAPTER – I: INTRODUCTION	1
1.1 Origin and history of the problem	1
1.2 Computer aided diagnosis (CAD)	2
1.3 Renal cell carcinoma	4
1.3.1 Clear Cell	5
1.3.2 Papillary	5
1.3.3 Chromophobe	6
1.3.4 Renal Oncocytoma	6
1.4 Cellular Staining	7
1.5 Quantum Dots (QD)	10
1.6 Clinical Decision Support Systems	12
1.7 Summary	13
CHAPTER – II: MOLECULAR PROFILING	14
2.1 Image quantification system for colorectal cancer risk assessment using quantum dots and molecular profiling	15
2.1.1 Image analysis and quantification system	18
2.1.2 Image processing module	19
2.1.3 Quantification and molecular profiling	20
2.2 Summary	21
CHAPTER – III: SEGMENTATION	24
3.1 Color segmentation.	24
3.1.1 K-means for color segmentation	26
3.1.2 Color segmentation using color maps	27
3.2 Region of Interest (ROI) segmentation	46
3.3 Nuclear cluster segmentation	47
3.3.1 Preprocessing	49
3.3.2 Concavity or notch detection	51
3.3.3 Cluster Segmentation	53
3.3.4 Cell-size computation	54
3.3.5 Cluster identification	54
3.3.6 Notch pairing based on distance threshold	54
3.3.7 Notch pairing using centroid	54
3.3.8 Ellipse fitting	55
3.3.9 Results	55
3.4 Summary	56
CHAPTER– IV: FEATURE EXTRACTION, SELECTION AND CLASSIFICATION	58
4.1 Knowledge based features	59
4.1.1 System Design	61
4.1.2 Tissue Samples and Image Collection	63
4.1.3 Image Processing	63
4.1.4 Feature Extraction	64
4.1.5 Classification	66
4.2 Morphological, textural and wavelets based features	68
4.2.1 Feature Extraction	68

4.2.2 Classification	73
4.3 Cellular features of elliptical models of segmented nuclei clusters	75
4.3.1 Methodology	78
4.3.2 Image acquisition & image color segmentation	79
4.3.3 Nuclei cluster segmentation and ellipse fitting	80
4.3.4 Feature extraction and selection	82
4.3.5 Classification	85
4.3.6 Results and discussion	88
CHAPTER – V : RCC GRADING BASED ON FUHRMAN NUCLEAR GRADE	90
5.1 Fuhrman Grading	90
5.2 Nuclear segmentation of high grade images	93
5.3 Feature extraction for Fuhrman grading	100
5.4 Scoring Nuclei	101
CHAPTER – VI: CDSS FOR FUHRMAN GRADING	106
6.1 Graphical user interface (GUI)	107
6.2 Time analysis	110
6.3 Grade predictions and decisions	111
6.4 Feedback and results	113
Summary	115
CONCLUSION	116
REFERENCES	118

## LIST OF TABLES

Table 2.1	A database table extract showing some of the section parameters .....	21
Table 3.1	Performance comparison of different visualization techniques.....	35
Table 3.3	Quantitative analysis of nuclear segmentation for four different types of RCC .....	56
Table 4.1	Feature extracted for Papillary (PAP), Clear Cell (CC), Chromophobe (CHR) and Renal Oncocytoma (ONC) with mean and standard deviations .....	65
Table 4.2	Features extracted from GLCM for clear cell (CC), chromophobe (CHR), renal oncocytoma (ONC) and papillary (PAP) with mean and standard deviations.....	70
Table 4.3	Statistical Features extracted for clear cell (CC), chromophobe (CHR), renal oncocytoma (ONC) and papillary (PAP) with mean and standard deviations.....	71
Table 4.4	Features extracted after DWT for clear cell (CC), chromophobe (CHR), renal oncocytoma (ONC) and papillary (PAP) with mean and standard deviations.....	72
Table 4.5	List of features selected for the best classification performance with mean and standard deviations. ....	73
Table 4.6	Confusion Matrix .....	87
Table 5.1	General guideline for Fuhrman nuclear grading. ....	91
Table 5.2	Features used for grading of RCC images .....	101
Table 6.1	Time analysis showing time required for different components for image examples of different sizes. ....	111

## LIST OF FIGURES

<b>Figure 1.1:</b>	Statistics related to new cancer cases in US (2010) with zoomed data for all sites and urinary cancers.....	2
<b>Figure 1.2:</b>	Photomicrograph of a clear cell RCC shows characteristic clear cytoplasm. ....	5
<b>Figure 1.3:</b>	Photomicrograph of a Papillary RCC showing characteristic white streaks.	6
<b>Figure 1.4:</b>	Photomicrograph of a chromophobe RCC shows characteristic peri-nuclear halos (arrows). ....	6
<b>Figure 1.5:</b>	Photomicrograph of a renal oncocytoma.....	7
<b>Figure 1.6:</b>	4-QD stained prostate tissue shown cellular structure of prostate glands..	11
<b>Figure 2.1</b>	Biomarkers of risk for colorectal cancer. (Left) Biomarker detection using traditional IHC. (Right) Biomarker detection using quantum dot IHC.....	16
<b>Figure 2.2</b>	Schematic representation of traditional and quantum dot IHC to detect biomarkers of risk for colorectal cancer.....	17
<b>Figure 2.3</b>	Components of integrated image analysis and quantification system .....	19
<b>Figure 2.4</b>	An application snapshot showing (a) hemi-crypt segmentation, (b) color separation, (c) sectioning, (d) GUI, (e) biomarker density distribution .....	20
<b>Figure 2.5</b>	A GUI with overlays showing QD stained prostate tissue with segmented gland, a single biomarker signature, pseudo-color biomarker signature representation, sample points for radial profiling and intensity profile of the biomarkers. ....	22
<b>Figure 2.6</b>	A GUI showing IHC stained prostate tissue, livewire segmentation of ROI and biomarker signatures.....	23
<b>Figure. 3.4</b>	(a) Histogram showing pixel count in input image binned according to the pixel color (b) 2D visualization of Histogram values (c) Space allocation for individual cell (d) Space allocation for complete map (e) Visualization map after radial sort.....	32
<b>Figure. 3.5</b>	(a) Original Image (b) Flattened Histogram.....	33
<b>Figure. 3.6</b>	Spiral search .....	33
<b>Figure 3.7</b>	Different stages of spiral spread scheme showing the evolution of map....	34
<b>Figure 3.8</b>	(a) Original Image (b) Random spread visualization of the original image .....	35
<b>Figure 3.16</b>	An example of cumulative cost and path matrix used for optimal path computation .....	47
<b>Figure 3.17</b>	An application snapshot showing Hemi-crypt segmentation using live-wire segmentation.....	47
<b>Figure. 3.19</b>	Pre-processing steps implemented for papillary tissue sample. (a) Input RGB image shown in gray scale, (b) binary mask of nuclei, (c), filled binary mask, (d) mask after noise removal (e) result after edge detection, (f) result after smoothing.....	49
<b>Figure. 3.18</b>	a) Overall flow-diagram for the method, b) Pre-processing steps .....	49
<b>Figure 3.20</b>	A synthetic cluster illustrating the method of calculating $\Theta$ , angle between adjacent normals.....	51
<b>Figure 3.21</b>	a) Graph illustrating angles between adjacent normals for different segment number along the edge of cluster in figure 3. Dotted line and complete line	



	represents high and low threshold respectively. b) Relation between z-component and segment number for the same cluster. Thin circles mark true concavities and thick circles mark the possible false concavities .....	52
<b>Figure 3.22</b>	Figure depicts cross product resultant direction in case of convex and concave contour locations; concavities are marked .....	53
<b>Figure 3.23</b>	Flow-chart for Segmentation.....	53
<b>Figure 3.24</b>	a) Input papillary tissue, b) result image after segmentation of image, green line mark the cell boundaries.....	55
<b>Figure 4.1</b>	Workflow of the proposed system showing data flowing between image acquisition, processing, feature extraction, classification, feedback and storage modules. ....	62
<b>Figure 4.2</b>	Scatter Plot showing distribution of Images for three co-occurrence features: Contrast, Correlation and Homogeneity (+, Papillary; o, Clear Cell; x, Chromphobe; *, Oncotytoma) .....	67
<b>Figure 4.3</b>	GLCM computation using 4-level grayscale images. (a) Representation of 4 level grayscale image (b) GLCM for highlighted elements in image (a) .....	69
<b>Figure 4.4</b>	(a) 4 level DWT image, (b) Level three approximation component, (c) Level two horizontal detail component, (d) 4-level grayscale ONC image. ....	72
<b>Figure 4.5:</b>	Scatter plot showing distribution of images for three features: GLCM based Energy, GLCM based Diagonal component representing cytoplasm area and Wavelet level 1 Diagonal detail GLCM component (2,2).....	74
<b>Figure 4.6:</b>	(Center) PAP image misclassified as ONC. (Left) another PAP image (Right) ONC image .....	74
<b>Figure 4.7:</b>	(Center) CC image misclassified as ONC. (Left) another CC image (Right) ONC image.....	74
<b>Figure 4.8</b>	Flowchart for the overall methodology.....	78
<b>Figure 4.9</b>	K-Means segmentation (Left to right): 1) original image; 2) gray level segmented image; 3) segmented pseudo color image; 4) Nuclei mask.....	80
<b>Figure 4.10</b>	(Left) Cluster segmentation by watershed method (Right) Cluster segmentation by notch detection and ellipse fitting method .....	80
<b>Figure 4.11</b>	Nuclei cluster segmentation for PA RCC tissue image. (a) RGB image, (b) Binary nuclei mask, (c) Individual nuclei marked on RGB tissue image using ellipse fitting. ....	82
<b>Figure 4.12</b>	(Left) CC input image. (Right) Mask of regions outside the nuclei.....	84
<b>Figure 4.13</b>	Statistics for top models. (a -top) Most significant features based on how frequently they showed up in top models (b - bottom left) Best K-value selection for KNN. (c -bottom right) Number of features used by top models .....	86
<b>Figure 4.14</b>	Distribution of normalized features for different subtypes .....	87
<b>Figure 4.15</b>	Correlation between cross validation and external validation results .....	87
<b>Figure 5.1</b>	RCC images of different grades showing variation in nuclear features. ....	92
<b>Figure 5.2</b>	Workflow for the proposed methodology of RCC grading .....	92
<b>Figure 5.3</b>	Higher grade image showing lightly stained nuclei due to chromatic activity .....	94
<b>Figure 5.4</b>	The parameter space used for CHT .....	95

<b>Figure 5.5</b>	A Circular Hough transform from the x,y-space (left) to the parameter space (right), this example is for a constant radius .....	96
<b>Figure 5.6</b>	<i>abcdef</i> (a) Synthetic image (b) Accumulator array for one specific radius (c) sum of all radii accumulator planes (d) Sum of accumulator after thresholding values (e) Binary mask of (d) (f) resultant circular models. ....	97
<b>Figure 5.7</b>	<i>a b c d</i> (a) Synthetic image (b) Sum of all radii accumulator planes (c) Resultant circular models. (d) Accumulator using radial lines .....	97
<b>Figure 5.8</b>	Flow chart for nuclear segmentation using gradient lines based approach	99
<b>Figure 5.9</b>	Results of gradient lines algorithm on synthetic data showing accumulator array, merging and splitting objects to detect nuclei.....	99
<b>Figure: 5.10</b>	<i>abcd</i> (a) Resultant image with Green circles showing detected nuclei and blue circles showing rejected objects(b) Input image(c) Nuclear mask of input image (d) Edges used for computation of Hough transform.....	100
<b>Figure: 5.11</b>	<i>abcde</i> (a) Nucleus image (b) Nuclear area (c) Nuclear area filled (d) Convex hull area of nuclei (e) Nuclear area unstained - red.....	101
<b>Figure 5.12</b>	Grade I –IV (left to right) representative images are shown with automatic selection of nuclei based on their size and eccentricity. The selected nuclei segmentation is also shown as overlay in bottom row. Holes in nuclear masks show nucleolus prominence. ....	103
<b>Figure 5.13</b>	Two views of clustering of different grade images based on the selected features i.e. nuclear convex area , nuclear circularity, nuclear solidity and unstained nuclear area. ....	104
<b>Figure 6.1</b>	GUI showing image loading, entry for basic parameters, 2D segmentation map and user instructions. ....	108
<b>Figure 6.2</b>	GUI showing zone marking by user for segmenation .....	108
<b>Figure 6.3</b>	GUI showing segmenation results in psuedo color.....	109
<b>Figure 6.4</b>	GUI showing top ranked nuclei(overlaid in green) based on their grading score(black). The holes in the segmentation masks show nucleolus prominence. Major statistics along with the predicted grade are shown in results panel. Feedback panel records final decision, system performance and specific comments. ....	109
<b>Figure 6.5</b>	Original and annotated images are shown which can be used for pathologists review .....	112
<b>Figure 6.7</b>	(a&b) Nuclei size for papillary carcinoma is larger than the one in corresponding clear cell reference. (a&c) For same size and nucleolus prominence image exhibiting pleomorphism shows higher grade. (d) Cell structure is clear because of prominent cell boundaries. (e) Nuclei are clustered together but cell boundaries aren't visible and bi-nucleation cannot be determined.....	114

## SUMMARY

This research focuses towards the development of a clinical decision support system (CDSS) based on cellular and tissue image analysis and classification system that improves consistency and facilitates the clinical decision making process. In a typical cancer examination, pathologists make diagnosis by manually reading morphological features in patient biopsy images, in which cancer biomarkers are highlighted by using different staining techniques. This process is subjected to pathologist's training and experience, especially when the same cancer has several subtypes (i.e. benign tumor subtype vs. malignant subtype) and the same cancer tissue biopsy contains heterogeneous morphologies in different locations. The variability in pathologist's manual reading may result in varying cancer diagnosis and treatment.

This Ph.D. research aims to reduce the subjectivity and variation existing in traditional histo-pathological reading of patient tissue biopsy slides through Computer-Aided Diagnosis (CAD). Using the CAD, quantitative molecular profiling of cancer biomarkers of stained biopsy images are obtained by extracting and analyzing texture and cellular structure features. In addition, cancer sub-type classification and a semi-automatic grade scoring (i.e. clinical decision making) for improved consistency over a large number of cancer subtype images can be performed. The CAD tools do have their own limitations and in certain cases the clinicians, however, prefer systems which are flexible and take into account their individuality when necessary by providing some control rather than fully automated system. Therefore, to be able to introduce CDSS in health care, we need to understand users' perspectives and preferences on the new information technology. This forms as the basis for this research where we target to present the quantitative information acquired through the image analysis, annotate the images and provide suitable visualization which can facilitate the process of decision making in a clinical setting.

# **CHAPTER - I**

## **INTRODUCTION**

This research focuses towards the development of a clinical decision support system (CDSS) based on cellular and tissue image analysis and classification system that improves consistency and facilitates the clinical decision making process. The topic has many facets and some introduction to all of these is deemed necessary towards the understanding of the complete problem. Keeping this in view, a brief review of all these areas have been covered in the succeeding sections of this chapter.

### **1.1 Origin and history of the problem**

Cancer is the second leading cause of death after heart disease in America. Although there has been a steady decrease in the incidence of death due to heart disease, no such trend can be observed in cancer. In addition to that, the risk of getting cancer is increasing due to major environmental, habitual and behavioral trends[1]. Some of the statistics published by American Cancer Society [2] in 2010, related to diagnosis and mortality (Figure 1.1) shows the immense scale of the problem.

Even with impressive strides made in treating and curing cancer, further improvement of survival rate relies heavily on the early diagnosis of cancer. Also, important to clinical success is to know the behavior of a particular cancer and its treatment, which depends on correct identification of cancer stage and/or subtype. It is therefore imperative in the fight against cancer that we not only diagnose cancer early, but also differentiate between its various subtypes quickly and accurately. To achieve relatively high differential accuracy, extensive training is usually required by a pathologist. Unfortunately, the current diagnostic paradigm of manual assessment of histology slides is slow and often irreproducible [3-5]. By leveraging the power of

computational systems, we can not only speed up the process of diagnosis, but also reduce the subjectivity in pathology.

Our research is targeted towards the design and development of a novel computer-aided diagnosis system which seeks interaction with expert users throughout the diagnosis process. With the system, users can lend their expertise to the validation of feature extraction and quantification, and they can also select from a list of features they deem most important and appropriate for the classification at hand. With such user interactivity and flexibility, this same tool is designed to be used by different pathologists with different cancer classification goals as long as the system has had sufficient training.

Estimated New Cancer Cases and Deaths by Sex for All Sites, US, 2010*						
	Estimated New Cases			Estimated Deaths		
	Both Sexes	Male	Female	Both Sexes	Male	Female
All Sites	1,529,560	789,620	739,940	569,490	299,200	270,290
Oral cavity & pharynx						
Tongue	43,140	21,370	21,770	16,800	8,770	18,030
Mouth	4,880	1,660	3,220	2,290	850	1,440
Other oral cavity						
Digestive system						
Esophagus	240,610	130,600	110,010	161,670	86,550	75,120
Stomach	12,220	10,110	2,110	3,800	2,870	730
Small intestine	222,530	116,750	105,770	157,300	86,220	71,080
Colon*	5,370	3,740	1,630	770	460	310
Rectum	2,650	1,530	1,120	1,440	890	550
Anal, anal canal, & anorectum						
Liver & intrahepatic bile duct	74,010	42,610	31,400	11,790	7,910	3,880
Gallbladder & other biliary	48,130	38,870	29,260	8,700	5,670	3,030
Pancreas	5,880	3,760	2,120	2,090	1,260	830
Other digestive organs						
Respiratory system						
Larynx	240,610	130,600	110,010	161,670	86,550	75,120
Lung & bronchus	12,220	10,110	2,110	3,800	2,870	730
Other respiratory organs	222,530	116,750	105,770	157,300	86,220	71,080
Breast & parts	5,370	3,740	1,630	770	460	310
Soft tissue (including heart)	2,650	1,530	1,120	1,440	890	550
Skin (excluding basal & squamous)	10,520	5,680	4,840	3,920	2,030	1,900
Melanoma*	74,010	42,610	31,400	11,790	7,910	3,880
Other non-melanoma skin	48,130	38,870	29,260	8,700	5,670	3,030
Breast	5,880	3,760	2,120	2,090	1,260	830
Genital system						
Uterine cervix	209,060	1,970	207,090	40,230	390	39,840
Uterine corpus	311,210	277,860	83,350	60,420	32,710	27,710
Ovary	12,200	12,200	0	6,210	6,210	0
Vagina	43,470	43,470	0	7,960	7,960	0
Vagina & other genital, female	21,880	21,880	0	13,860	13,860	0
Prostate	3,900	3,900	0	920	920	0
Penis & other genital, male	2,500	2,500	0	780	780	0
Urinary system						
Urinary bladder	217,730	217,730	0	32,050	32,050	0
Kidney & renal pelvis	8,480	8,480	0	350	350	0
Ureter & other urinary organs	1,250	1,250	0	310	310	0
Urinary system	131,260	89,620	41,640	28,550	19,110	9,440
Urinary bladder	70,530	70,530	0	14,680	14,680	0
Kidney & renal pelvis	58,240	35,370	22,870	13,040	8,210	4,830
Ureter & other urinary organs	2,490	1,490	1,000	830	490	340
Eye & orbit						
Brain & other nervous system						
Endocrine system						
Thyroid						
Other endocrine						
Lymphoma						
Hodgkin lymphoma						
Non-Hodgkin lymphoma						
Myeloma						
Leukemia						
Acute lymphocytic leukemia						
Chronic lymphocytic leukemia						
Acute myeloid leukemia						
Chronic myeloid leukemia						
Other leukemia*						
Other & unspecified primary sites*						

**Figure 1.1:** Statistics related to new cancer cases in US (2010) with zoomed data for all sites and urinary cancers

## 1.2 Computer aided diagnosis (CAD)

In the past, much effort has been devoted to development of automated diagnosis systems for various maladies. Although the application of computational methods represents a significant step forward in the fight against cancer and its early detection, the task has been anything but trivial. Traditionally, most automated cancer diagnosis

research has been on the problem of differentiating between cancerous and normal tissue images [6, 7].

Since pathologists use deviations in cellular structure as a means to make a diagnosis, many previous methods have used the statistical variation of various image properties to help make a diagnosis. The use of morphological features, for example, was reported by Jiang et al. in their study of breast cancer classification and by Roula et al. for the grading of prostate cancer [8, 9].

The diagnosis system developed by Diamond et al. used a combination of structural and textural features to achieve an accuracy of 79.3% for the classification of prostatic neoplasia [10]. Esgiar et al. studied the classification of colonic mucosa using six different textural features and optical density and reported an overall accuracy of 90.2 % [11].

The choice of features in these studies was generally motivated by the hypothesis that the human eye uses these features for discrimination and so should automated systems. Recent attempts at segmentation have moved beyond this paradigm and have instead included less intuitive features such as fractal dimension which are not easily detectable or describable by humans. In a follow-up study, Esgiar et al. reported an increase in accuracy of their system when fractal analysis was employed along with previously identified textural features. Furthermore, they suggested that knowledge incorporation is needed for further increase in accuracy [12]. Hamilton et al. used such knowledge-guided segmentation to calculate features like the co-occurrence matrix and optical density to achieve 83% correct classification of colorectal dysplasia [13].

We agree with Esgair's assertion, and contend that classification accuracy can be increased by incorporating knowledge from an expert pathologist at every step of the system: image processing, feature extraction, and classification. Also, involving the user in the decision making process and allowing the user to bias the system will lead to making a more accurate prediction. The belief is that having user interaction through

every step of the process helps to encompass the vast heterogeneity found in tissue imaging data, thus making the system more robust to intra-class variation.

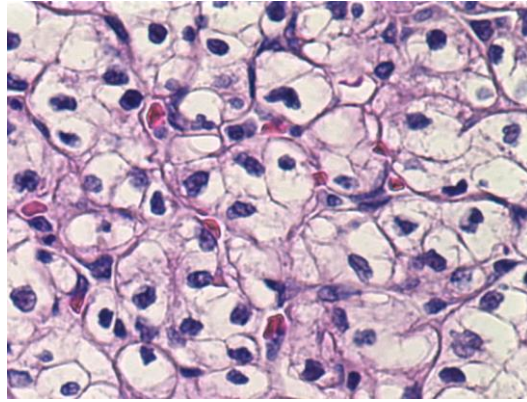
### **1.3 Renal cell carcinoma**

We chose renal cell carcinoma (RCC) as a case study for the development of this tool primarily because little research has been done for the automated classification of RCC. Moreover, this problem is more complicated than simple normal versus cancerous tissue classification as RCC has four commonly observed clinical subtypes. Moreover, early stage kidney cancer usually has no symptoms[2]. Symptoms that may develop as the tumor progresses include blood in the urine, a pain or lump in the lower back or abdomen, fatigue, weight loss, fever, or swelling in the legs and ankles.

RCC is the most common form of kidney cancer arising from the renal tubule in adults [2, 14]. An estimated 58,240 new cases of kidney (renal) cancer were diagnosed in 2010 and an estimated 13,040 deaths from kidney cancer occurred in 2010. RCC begins small and grows larger over time, like many other cancers. RCC usually grows as a single mass. Sometimes, a kidney may contain more than one tumor or tumors may be found in both kidneys at the same time. There are several sub-types of renal cell cancer and the prognosis and treatment can depend on what sub-type you have. More than 90% of clinically significant lesions can be diagnosed as one of the common subtypes of renal tumor: clear cell RCC (CC), papillary RCC (PAP), chromophobe RCC (CHR), and renal oncocytoma (ONC). We will explore some of the features of these subtypes as these will be used repeatedly in the succeeding chapters of this document.

### 1.3.1 Clear Cell

This is by far the most common sub-type of the RCC [15] and more than 75% of the renal cell cancers belong to this subtype. CC predominantly consists of cells with



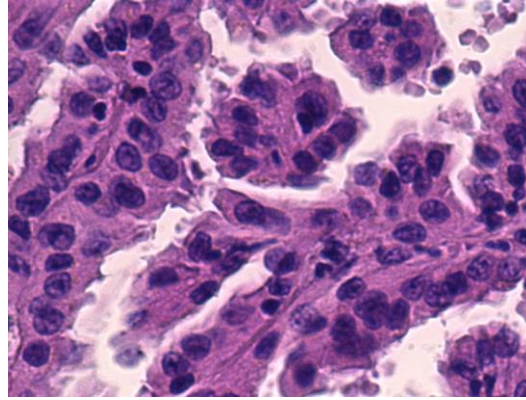
clear cytoplasm.

**Figure 1.2:** Photomicrograph of a clear cell RCC shows characteristic clear cytoplasm.

### 1.3.2 Papillary

Papillary renal cell carcinoma is the second most common sub-type after clear cell. About 12 % of the renal tumors belong to this subtype[15] . These cancers form little finger-like projections (called papillae) in most of the tumor. Some doctors call these cancers *chromophilic* because the cells take in certain dyes used so the tissue can be seen under the microscope, and look pink.

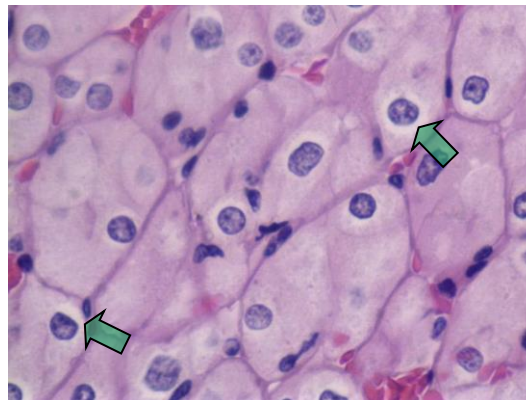




**Figure 1.3:** Photomicrograph of a Papillary RCC showing characteristic white streaks.

### 1.3.3 Chromophobe

Chromophobe RCC occurring approx. 4% amongst renal cancer[15], tends to metastasize to the liver more than clear cell RCC. Chromophobe RCC is histopathologically characterized by large polygonal cells with prominent cell membranes (Figure 1.4). In contradistinction to clear cell RCC, the tumor blood vessels are thick walled and eccentrically hyalinized

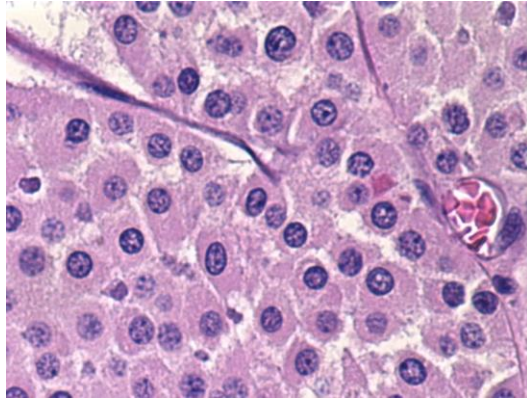


**Figure 1.4:** Photomicrograph of a chromophobe RCC shows characteristic peri-nuclear halos (arrows).

### 1.3.4 Renal Oncocytoma

Renal Oncocytoma is actually a benign tumor. Renal Oncocytoma occurs at about 5 to 10% [15] of the rate of kidney cancer. Renal oncocytoma can occasionally be confused with chromophobe RCC or the granular variant of clear cell RCC, although

most cases are easily diagnosed. Although benign, oncocytoma occasionally co-exists with cancer which may be present within or near the oncocytoma.



**Figure 1.5:** Photomicrograph of a renal oncocytoma.

Renal tumor subtypes exhibit several common morphological characteristics, making diagnosis difficult and subjective in many cases. Histopathologic classification is critical for the treatment of RCC because each of its subtypes is associated with a distinct clinical behavior. Development of a diagnosis technique with a quantitative approach to renal tumor classification is therefore critical and very much needed.

Expert knowledge of RCC features was incorporated into our system by letting a pathologist select the features most relevant to him for the diagnosis of RCC. This was coupled with the prior knowledge about the presence and/or absence of specific histological features and structures (Red blood cells, blood vessels, lipid structures, papillary bodies) in various subtypes of RCC.

#### **1.4 Cellular Staining**

Staining [16] is an auxiliary technique used in microscopy to enhance contrast in the microscopic image. Stains and dyes are frequently used in biology and medicine to highlight structures in biological tissues for viewing, often with the aid of different microscopes. Stains may be used to define and examine bulk tissues (highlighting, for example, muscle fibers or connective tissue), cell populations (classifying different blood cells, for instance), or organelles within individual cells.

In biochemistry it involves adding a class-specific (DNA, proteins, lipids, carbohydrates) dye to a substrate to qualify or quantify the presence of a specific compound. Staining and fluorescent tagging can serve similar purposes. Biological staining is also used to mark cells in flow cytometry, and to flag proteins or nucleic acids in gel electrophoresis.

Staining is not limited to biological materials, it can also be used to study the morphology of other materials for example the lamellar structures of semi-crystalline polymers or the domain structures of block copolymers.

Cell staining techniques [17] and preparation depend on the type of stain and analysis used. One or more of the following procedures may be required to prepare a sample:

**Permeabilization** - treatment of cells, generally with a mild surfactant, which dissolves cell membranes in order to allow larger dye molecules to enter inside the cell.

**Fixation** - serves to "fix" or preserve cell or tissue morphology through the preparation process. This process may involve several steps, but most fixation procedures involve adding a chemical fixative that creates chemical bonds between proteins to increase their rigidity. Common fixatives include formaldehyde, ethanol, methanol, and/or picric acid.

**Mounting** - involves attaching samples to a glass microscope slide for observation and analysis. Cells may either be grown directly to the slide or loose cells can be applied to a slide using a sterile technique. Thin sections (slices) of material such as tissue may also be applied to a microscope slide for observation.

**Staining** - application of stain to a sample to color cells, tissues, components, or metabolic processes. This process may involve immersing the sample (before or after fixation or mounting) in a dye solution and then rinsing and observing the sample under a microscope. Some dyes require the use of a mordant, which is a

chemical compound that reacts with the stain to form an insoluble, colored precipitate. The mordanted stain will remain on/in the sample when excess dye solution is washed away.

There are several types of staining media; each can be used for a different purpose. Commonly used stains and how they work are listed below. All these stains may be used on fixed, or non-living, cells and those that can be used on living cells are noted.

**Bismarck Brown** - colors acid mucins, a type of protein, yellow and may be used to stain live cells

**Carmin** - colors glycogen, or animal starch, red

**Coomassie blue** - stains proteins a brilliant blue, and is often used in gel electrophoresis

**Crystal violet** - stains cell walls purple when combined with a mordant. This stain is used in Gram staining

**DAPI** - a fluorescent nuclear stain that is excited by ultraviolet light, showing blue fluorescence when bound to DNA. DAPI can be used in living or fixed cells

**Eosin** - a counterstain to haematoxylin, this stain colors red blood cells, cytoplasmic material, cell membranes, and extracellular structures pink or red.

**Ethidium bromide** - this stain colors unhealthy cells in the final stages of apoptosis, or deliberate cell death, fluorescent red-orange.

**Fuchsin** - this stain is used to stain collagen, smooth muscle, or mitochondria.

**Hematoxylin** - a nuclear stain that, with a mordant, stains nuclei blue-violet or brown.

**Hoechst stains** - two types of fluorescent stains, 33258 and 33342, these are used to stain DNA in living cells.

**Iodine** - used as a starch indicator. When in solution, starch and iodine turn a dark blue color.

**Malachite green** - a blue-green counterstain to safranin in Gimenez staining for bacteria. This stain can also be used to stain spores.

**Methylene blue** - stains animal cells to make nuclei more visible.

**Neutral/Toluylene red** - stains nuclei red and may be used on living cells.

**Nile blue** - stains nuclei blue and may be used on living cells.

**Nile red/Nile blue oxazone** - this stain is made by boiling Nile blue with sulfuric acid, which creates a mix of Nile red and Nile blue. The red accumulates in intracellular lipid globules, staining them red. This stain may be used on living cells.

**Osmium tetroxide** - used in optical microscopy to stain lipids black.

**Rhodamine** - a protein-specific fluorescent stain used in fluorescence microscopy.

**Safranin** - a nuclear stain used as a counterstain or to color collagen yellow.

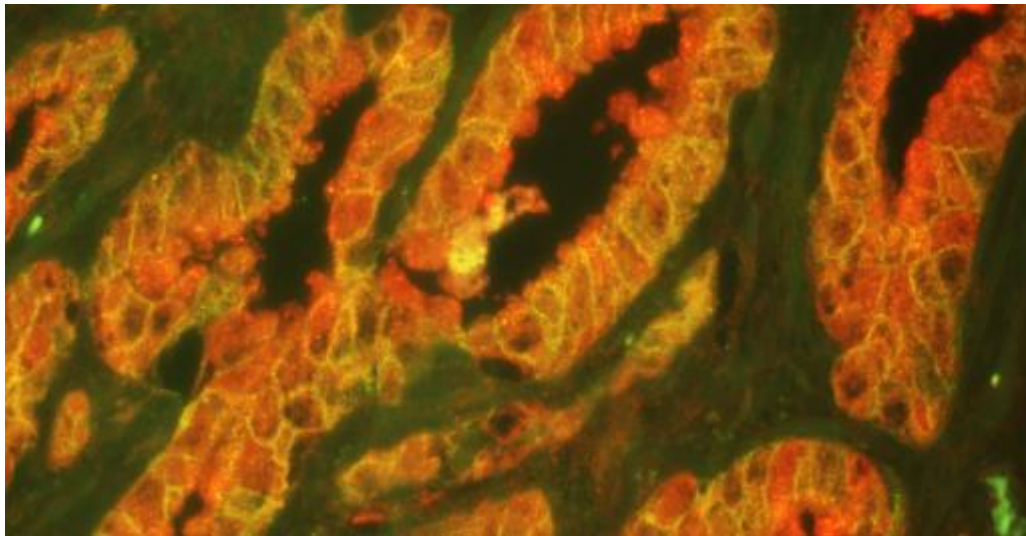
After staining cells and preparing slides, they may be stored in the dark and possibly refrigerated to preserve the stained slide, and then observed with a microscope.

## 1.5 Quantum Dots (QD)

QDs are tiny light-emitting particles on the nanometer scale, and are emerging as a new class of fluorescent labels for biology and medicine [18, 19]. In comparison with organic dyes and fluorescent proteins, QDs have unique optical and electronic properties such as size-tunable light emission, superior signal brightness, resistance to photo-bleaching and simultaneous excitation of multiple fluorescence colors. These properties are most promising for improving the sensitivity and multiplexing capabilities of molecular histopathology and disease diagnosis. Recent advances have led to highly bright and stable QD probes that are well suited for profiling genetic and protein biomarkers in intact cells and clinical tissue specimens[20]. In contrast to in vivo imaging

applications where the potential toxicity of cadmium-containing QDs is a major concern, immunohistological staining is performed on in vitro clinical patient samples. As a result, the use of multicolor QD probes in immunohistochemistry (IHC) is likely one of the most important and clinically relevant applications in the near term.

In recent years, several groups have used QD probes for fluorescence immunostaining of fixed cells and tissue specimens[21]. However, medical applications of QD-based immunohistochemistry have not achieved widespread adaptation or significant clinical success. A major problem is the lack of robust protocols and experimental procedures to define the key factors and steps involved in QD immunohistochemical staining and data analysis. In particular, there are no consensus on methods for QD– antibody (QD–Ab) bioconjugation, tissue specimen preparation, multicolor QD staining, image processing and data quantification. We collaborated in one such effort [22] for development of antibody- conjugated QDs for multiplexed and quantitative (or semi-quantitative) IHC, and five-color molecular profiling on formalin-fixed and paraffin-embedded (FFPE) clinical tissue specimens was have achieved. An example of QD stained prostate tissue is shown in Figure 1.6



**Figure 1.6:** 4-QD stained prostate tissue shown cellular structure of prostate glands.

## 1.6 Clinical Decision Support Systems

CDSS are computer systems designed to facilitate clinician decision making process. The systems try to make individual specific decisions utilizing additional data and information available through their personal records and other sources. CDSS generally have three major components i.e. the knowledge base, the reasoning engine and the user interface [23]. The knowledge base contains the rules and associations of compiled data. The inference engine combines the rules from the knowledge base with the patient's data.

CDSSs in healthcare [24] have met with varying amounts of success in different domains including image analysis of radiology[25] [26] and histology images[27] [28]. Most successful areas include pharmacy and billing sectors while core areas are still the main focus of the community. According to Ken Congdon [29], CDSS is one of the top 10 IT trends in the year 2012(Table1.1) . At present focus of CDSS [30] is on following areas:

- Alerts and Reminders
- Diagnostic Assistance
- Prescription Decision Support
- Information Retrieval
- Image Recognition and Interpretation
- Therapy Critiquing and Planning

Table 1.1 Top 10 IT Trends for 2012

Trend Number	Category	Top Priority	Priority	Important	Somewhat Important	Not Important At This Time	Average Ranking
1.	EHR Adoption & Meaningful Use	55.5%	16.8%	11.5%	4.7%	11.5%	2.00
2.	HIPAA 5010 Compliance	39.6%	31.3%	14.1%	5.2%	9.8%	2.15
3.	ICD-10 Compliance	36.5%	27.6%	18.2%	6.8%	10.9%	2.28
4.	PHI Security	31.2%	23.8%	28.0%	8.5%	8.5%	2.39
5.	Clinical Decision Support	27.5%	24.9%	30.2%	7.9%	9.5%	2.47
6.	e-Prescribing	28.2%	28.2%	21.0%	8.0%	14.4%	2.52
7.	Wireless/Mobile Computing	18.8%	33.9%	32.3%	6.8%	8.3%	2.52
8.	Document Imaging/Management	16.5%	37.8%	25.5%	11.2%	9.0%	2.59
9.	HIE	20.9%	35.1%	22.0%	6.8%	15.2%	2.60
10.	Telehealth	19.5%	30.5%	28.4%	12.6%	8.9%	2.61

The investigation, analysis and interpretation of the pathological imaging data mainly depends on the pathologist's knowledge, experience and his subjective view about the

data. Although Computer Aided Diagnostic (CAD) tools [31] can help reduce this subjectivity and the inter-user and intra-user variability and gained some acceptance among clinicians [3, 4], sometimes the clinicians, prefer systems which are flexible and take into account their individuality when necessary by providing some control rather than fully automated system [32] [33]. Studies [34, 35] have shown improvement in practitioner's performance and patient's outcome for CDSS which account for practitioner's perspective and integrate it in the workflow. Therefore, to be able to introduce CDSS in health care [36] [37] [38], we need to understand users' perspectives and preferences on the new information technology. This forms as the basis for this research where we target to present the quantitative information acquired through the image analysis, annotate the images and provide suitable visualization which can facilitate the process of decision making in a clinical setting.

## **1.7 Summary**

In this chapter, we covered the origin and background of the problem. A brief review of CAD systems used to solve similar problems was carried out followed by some introduction to RCC basic staining techniques and CDSS trends which will provide us with enough background information to understand the work presented in the subsequent chapters.



## **CHAPTER - II**

### **MOLECULAR PROFILING**

Traditionally, tumors have been categorized on the basis of histology. The clinical behavior of human cancer has been predicted using its microscopic appearance. This has been a useful approach because cancer has hundreds of shapes and structures under the microscope. The staining pattern of cancer cells viewed under the microscope is insufficient to reflect the complicated underlying molecular events [39]. For example, prostate carcinoma arising in two patients may look virtually identical under the microscope, but each patient may have a different clinical outcome. This traditional classification scheme is also limited by a number of factors[40]. First, it relies on a subjective review of the tissue that is dependent on the knowledge and experience of a pathologist, and therefore may not be reproducible. The classification is discrete, rather than continuous, meaning that patients are classified into broad treatment groups (e.g., low, medium, or high probability of recurrence) with limited ability to determine the individual recurrence risk. In addition, current pathology reports either lack or offer very little information regarding the potential drug treatment regime to which a cancer will respond. While current pathology does help determine treatment that leads to better outcomes, tumors with identical pathology may have different origins and respond differently to treatment. Consequently there has been a persistent need to find some way to accurately subcategorize, and understand, the biological diversity of cancer.

Classification of cancerous tissue based on its molecular profile overcomes these limitations. A molecular profile determines the level of gene expression within the cancer by hybridizing the cellular RNA with known genes. Currently this is done using microarray technology to provide information on thousands of genes simultaneously. Once the gene expression pattern is determined, this information is compared to the

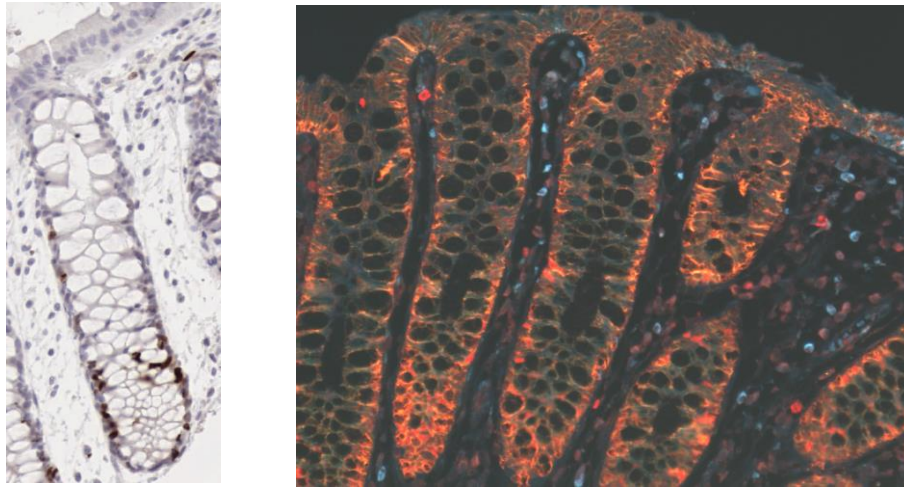
expression profiles of cancers with known outcomes using a predetermined algorithm. The algorithm then places the cancer into an outcome class based on similar gene expression patterns, or it will return a survival probability. However, the use of these tests for clinical decision making presents many challenges to overcome[41]. Assay development and data analysis in this field have been largely exploratory, and leave numerous possibilities for the introduction of bias. Optimal incorporation into clinical practice is not straightforward. Finally, cost-effectiveness is difficult to appreciate until these other challenges are addressed. Overall, molecular profiling is a fascinating and promising technology, but its incorporation into clinical decision making requires careful planning and robust evidence.

The information available under the microscope also has so much detail that mere pathological observation may not reveal the true information content. Moving away from the classic definition of molecular profiling, using advanced image analysis of the biomarker distributions and appropriate experimental models, the ultimate goal to move beyond correlation and classification to achieve new insights into disease mechanisms and treatment targets may be achieved. Keeping this in view, we carried out a case study[42] for the colorectal cancer where no known biomarkers of risk are available which can be used for predicting and preventing the disease.

## **2.1 Image quantification system for colorectal cancer risk assessment using quantum dots and molecular profiling**

Based on new knowledge of the molecular basis of colorectal cancer, we developed and validated a panel of biomarkers of risk that can be measured in rectal biopsies. The goal of this work is to develop an integrated detection and image analysis quantification system for measuring and applying these biomarkers in clinical research and care. More importantly, the new system can process biopsy images from both traditional and bio-nanotechnology quantum dot-based IHC, and through a combination

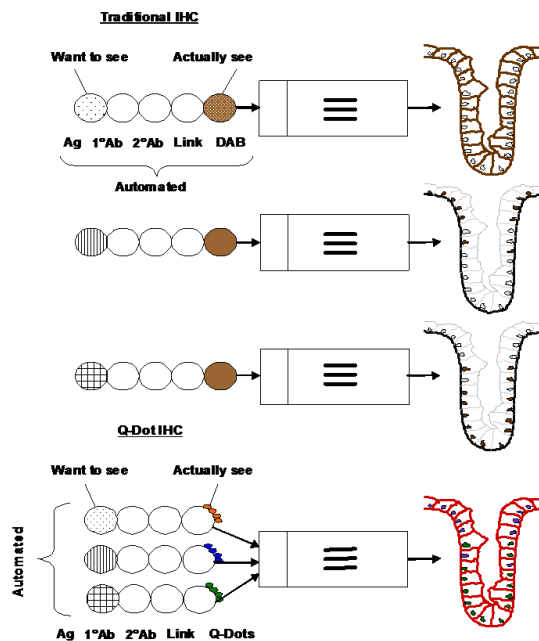
of novel and automated image analysis and quantification algorithms, it will significantly reduce processing time by detecting multiple biomarkers simultaneously on the same histologic sections. Clinical application of this novel process of detecting and quantifying biomarkers, coupled with decision support from the analysis of a biomarker quantification database, is expected to open new frontiers in the field of colorectal cancer prognosis and treatment.



**Figure 2.1** Biomarkers of risk for colorectal cancer. (Left) Biomarker detection using traditional IHC. (Right) Biomarker detection using quantum dot IHC

From rectal biopsies, we are interested in measuring the intensity and distributions of labeled antigens in the colon crypts. Colon crypts are test tube-like inversions of the inner colon lining. They are sites where colon polyps develop. Our existing system [43] uses immunohistochemical methods to detect the biomarkers in histologic slides (Figure 1) of normal colon tissue. In outline, immunohistochemistry (IHC) is a procedure in which an antigen (e.g., a protein biomarker) in a tissue is identified in a series of steps (Figure 2), including the application of a primary antibody to the antigen, linking the primary antibody to a secondary antibody that has an attachment site for a chemical linking agent, adding the linking agent that has an attachment site for a chromogen (e.g., DAB), and then adding the chromogen (one can also apply a counterstain at this point). Only one antigen can be detected on a given slide using this method. To detect six

different biomarkers, a set of five slides must be processed for each biomarker separately, for a total of 30 slides (6 biomarkers x 5 slides per set); each set of slides must also be analyzed separately. We also carried out experiments (Figure 2) using specially coated nanocrystals, called quantum dots (Q-Dots), instead of the chromogen. Q-Dots have the property of being easily excited to emit light in a very narrow spectrum. Q-Dots of slightly different sizes emit different, non-overlapping spectra. Q-Dots can be conjugated to the usual linking agents used in traditional IHC (Figure 2). This means that we can link Q-Dot-linking agent complexes with different size Q-Dots to different antibodies and thus to different antigens, thereby allowing detection of multiple biomarkers on the same slides. Also, in contrast to immunofluorescent dyes, the light emissions from quantum dots last months rather than just a few minutes, thus making analysis feasible in population- or clinical-based studies. With our new nanotechnology-based methods all six biomarkers can be detected at the same time (“multiplexed”) on one set of five slides (Figure 2).

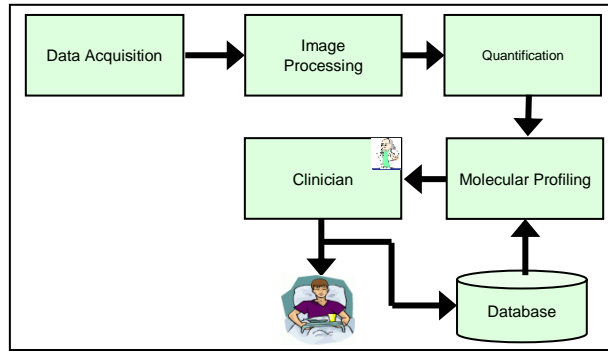


**Figure 2.2** Schematic representation of traditional and quantum dot IHC to detect biomarkers of risk for colorectal cancer

In our current system [43], candidate crypts are selected and traced by a trained technician. Biopsies are selected for scoring if they contain complete cross-sections of colon crypts and adequate staining is present in the slide. Then the technician draws a tracing encompassing one side of the crypt wall, termed a ‘hemi-crypt’. Colon cells within the tracing are segmented and measurements of the colon cells are recorded. This system is good enough for a research setting but certain limitations like capability to analyze only gray scale images, manual crypt tracing, and non-integrated systems components make it impractical for clinical use.

### **2.1.1 Image analysis and quantification system**

We have developed a complete integrated solution for image processing, quantification, and analysis of biomarkers for subsequent use in a clinical trial. The system (Figure 3) is capable of handling imaging data from both IHC and Q-Dot-based imaging. Imaging data are acquired by scanning multiple slides using an ultra-fast slide scanning System (T3 - Aperio Technologies), and selected colon crypts are processed using the image processing module. The module carries out semiautomatic hemi-crypt segmentation, color-based tissue classification, and hemi-crypt sectioning (discussed in detail in the following sections). Various biomarker features obtained through the image processing module are quantified and stored in a database. The obtained feature set is then analyzed and correlated to the existing data for predicting the results. The results will be validated by a clinician and feedback will be used to train the system for subsequent prediction. It is expected that the iterative process will eventually stabilize and the system’s prediction accuracy will increase to a limit where it can be used without a clinician’s validation.



**Figure 2.3** Components of integrated image analysis and quantification system

### 2.1.2 Image processing module

This module, developed in Matlab, is a collection of various image processing tools and is capable of processing image files of selected hemi-crypts from the Slide Scanner. The user interacts with the system through a graphical user interface (GUI) for performing various tasks (explained below) easily and efficiently.

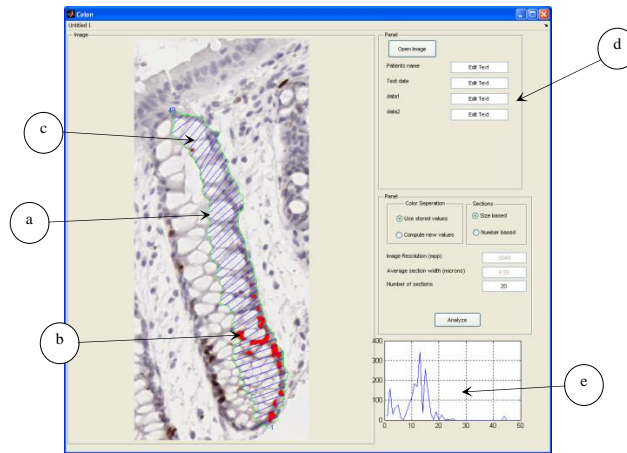
#### Semi-Automatic Hemi-crypt Segmentation and Color Segmentation

The first step in processing the images is to mark the region of interest, i.e., a hemi-crypt (one half of the symmetric colon crypt). This is accomplished by using our semi-automatic segmentation tool which has been discussed in detail in section (3.2) followed by the color segmentation classification of each pixel to belong to one of the biomarkers or to the background. This is achieved by using K-means clustering as described in section (3.1.1).

#### Hemi-crypt Sectioning

All of our colon cancer risk biomarkers are expressed in terms of density gradients along the lengths of colon crypts. As the shape and spatial orientation of the crypts can vary significantly from sample to sample, it poses a significant challenge to automatically determine the biomarker distribution and correlate this information with information obtained from other samples. We found the solution to this problem by dividing the hemi-crypt into subsections of uniform width. Two ends of the crypt are

marked by the user while tracing the hemi-crypt boundary during the segmentation process. Two user selectable modes are available to perform sectioning. These modes provide the option to section the hemi-crypt into a fixed number of slices or to partition it based on a user provided thickness parameter. Hemi-crypt sectioning is performed by computing the center line for the crypt, catering for all types of perimeter and shape variations, and ensuring no overlapping regions. The sectioning information is not only useful for determining the biomarker's density distribution, but also assists in maintaining standardization between samples for correlation purposes.



**Figure 2.4** An application snapshot showing (a) hemi-crypt segmentation, (b) color separation, (c) sectioning, (d) GUI, (e) biomarker density distribution

### 2.1.3 Quantification and molecular profiling

Because the amount of staining of a biomarker is proportional to the amount of biomarker in the tissue, and because the optical density of the staining is proportional to the amount of staining, we can quantify the amount of biomarker in the tissue using optical density measurements. The biomarkers thus detected during the image processing steps are converted into numeric format during the quantification process. Key parameters like biomarker intensities, background intensities, relationship of biomarker area to that of background area within the section, length of crypts etc. are of vital

importance. These parameters pertaining to the complete hemi-crypt as well as each subsection are stored in the database for subsequent analysis and classification.

Table 2.1 A database table extract, showing some of the section parameters

Hemicypt Number	Section Number	Section Area	Biomarker Area	Biomarker Avg Red intensity	Background Avg Red Intensity	Area Ratio
01	1	258.375	31.125	0.3539	0.6394	0.1205
01	2	779.500	152.250	0.2866	0.6858	0.1953
01	3	828.875	52.125	0.3246	0.7712	0.0629
01	4	761.125	96.625	0.3355	0.8373	0.127
01	5	851.000	9.875	0.4453	0.8094	0.0116

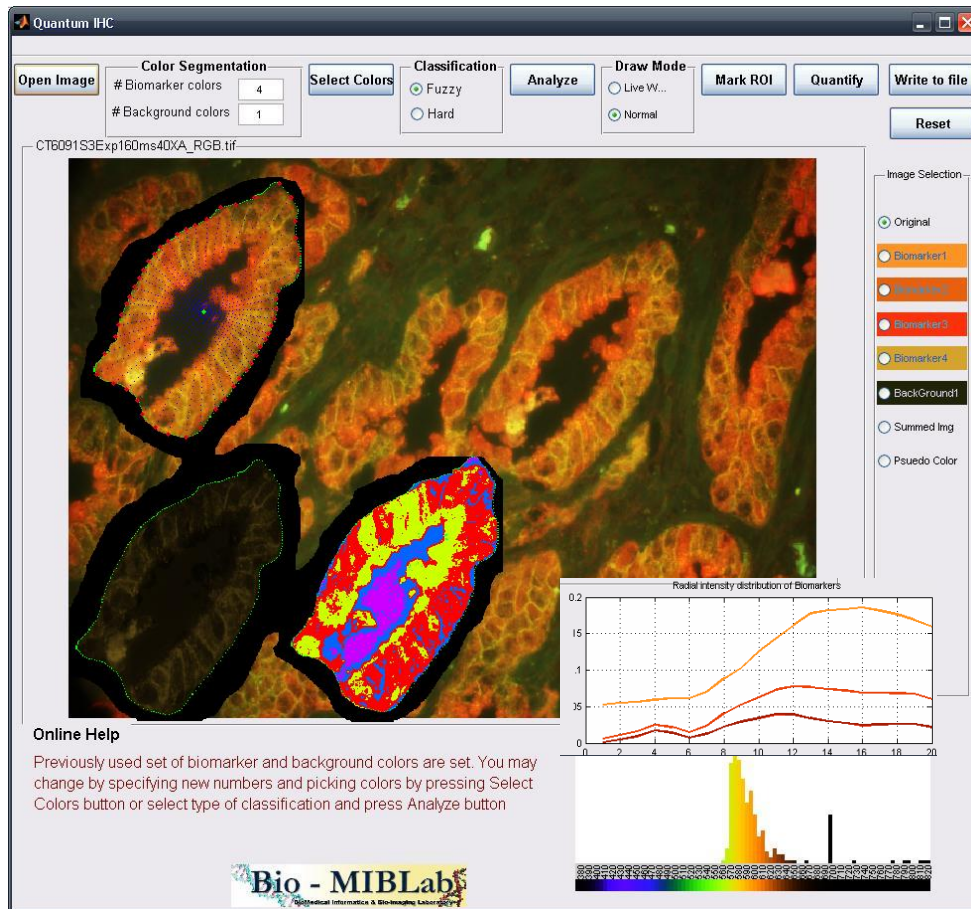
The analysis of these parameters will help us understand the molecular anatomy of normal cells and cells in various stages of progression to cancer. This will be achieved by correlating the biomarker density profiles (Figure 5e) of the patients with different datasets (already validated and stored in the database) to generate a prediction score for each sample and provide this score as a final output of our system.

## 2.2 Summary

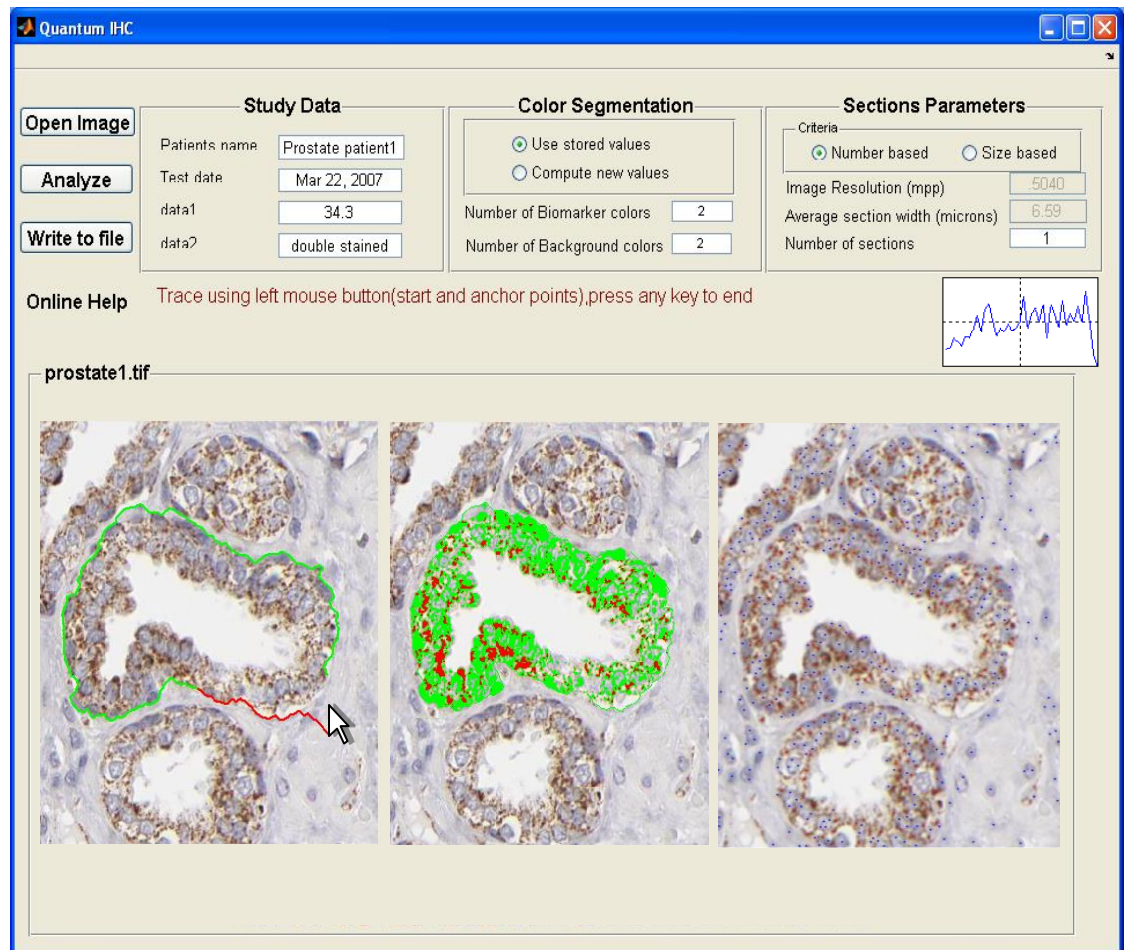
We have demonstrated that a valid panel of biomarkers of risk for colorectal cancer can be detected in histologic slides of normal colon tissue using traditional IHC as well as multiplexed, quantum dot IHC. Clinical application of these biomarkers requires a time efficient, consistent and accurate detection and processing system. These requirements are amply supported by our integrated image processing, quantification and analysis system. Using this novel implementation detection scheme will not only reduce the processing time and user effort significantly but will also add accuracy and consistency to the results. It is estimated that a single user will be able to process a set of slides in less than 5 minutes (excluding slide preparation time). These tools, using the semi-automated segmentation, will also reduce inter-user variability significantly thereby producing consistent results. In addition our color based biomarker detection is more accurate and easy to validate than our previous gray scale based detection system [43].



We extended this work to analyze prostate cancer specimens using both IHC and QD staining. We were able to successfully quantify different biomarkers and analyze their distributions and molecular profiles. Figure 2.5 and Figure 2.6 shows couple of examples for the GUI used for the analysis of prostate cancer. With all these improvements, our proposed system shows a great promise for clinically valid and practical methods of assessing and managing risk for cancer.



**Figure 2.5** A GUI with overlays showing QD stained prostate tissue with segmented gland, a single biomarker signature, pseudo-color biomarker signature representation, sample points for radial profiling and intensity profile of the biomarkers.



**Figure 2.6** A GUI showing IHC stained prostate tissue, livewire segmentation of ROI and biomarker signatures

## **CHAPTER - III**

### **SEGMENTATION**

Analysis of an image begins with a segmentation process, which differentiates meaningful regions of interest from the background. In our case, we are attempting to identify regions which most likely correspond to cells or nuclei. This step is critical[44] in that its outcome serves as the basis for all subsequent analyses, such as the extraction of shape features, and ultimately the interpretation of cell behavior and diagnosis.

Depending on segmentation requirements different techniques are adopted which are most suited to achieve the desired results. During our study, we employ segmentation for three different requirements i.e. color segmentation, ROI segmentation and nuclear cluster segmentation. Color segmentation is used to identify different regions based on their staining color. ROI segmentation, as the name implies, segments regions where analysis is required to be performed. The third requirement, nuclear cluster segmentation, works on already color segmented nuclear masks and uses shape based features to split nuclear clusters. We will discuss the segmentation techniques for all these focus areas in detail in following paragraphs.

#### **3.1 Color segmentation.**

Color segmentation is an important topic in biomedical image analysis because many biological discoveries and medical decisions are based on color staining of samples. Among many different aspects, the most relevant to our study are:

- a. Existing color segmentation applications in biomedicine (e.g. skin lesions etc.)
- b. Color representation and perception
- c. Interactivity and usability

Due to the complexity and heterogeneity of biomedical images, generic solutions are quite uncommon. A survey of biomedical image processing reveals a number of application-specific color segmentation techniques. One automatic color segmentation algorithm [45] is specifically designed for skin tumor feature identification. In that work, the segmentation depends on a database of feature information created by a dermatologist using specific software. Their methodology works well in feature based detection but can't be applied generically to a wide variety of images. Another work related to skin color segmentation [46] is based on a mixture of Gaussians. Satisfactory results were obtained for human skin color under different illumination. However, the method has not been tested for segmentation of images with appreciable color variation. Another segmentation algorithm, particularly related to skin color detection [47], but without any user interaction, depends on variable parameters and gives results which may not meet the user's expectation. A linear color segmentation technique [47] used to segment an image at material boundaries may suffer from spatial color heterogeneities which arise in supposedly uniform pigmented objects. Another work pertaining to color segmentation of pigmented skin lesions [48] is based on two-dimensional histogram analysis and a fuzzy K-means clustering technique. There, median filtering and morphological operations were used to smooth the border before segmentation. This approach might be useful for the particular problem addressed, (detection of pigmented skin lesions), but it does not address the major problem of color segmentation in biological images which may vary in color, intensity, imaging modality etc.

Images with sharp color distinctions can be easily segmented by using one of several segmentation algorithms available [49]. Current digital cameras make it possible to capture high resolution color images in clinical applications. However, natural and biological stained color images lack high contrast discontinuities. In addition, there are variations among staining colors or light conditions. Thus, using computers to conduct color quantification becomes increasingly important in clinical diagnosis, and emerges as

a new field called Computer Assisted Diagnosis (CAD). One of the challenges in CAD is how to perform precise color segmentation.

Gray-scale image segmentation use discontinuity-based techniques [50-52] to partition an image by detecting isolated points, lines and edges caused by sudden changes in gray levels. Homogeneity-based methods perform thresholding, clustering, region growing, and region splitting and merging [52]. Color is a feature of an object's surface. Distinctive colors form peaks in color histograms. Because color images are usually measured and represented by three color components: red green and blue, (R,G,B), most color segmentation approaches just extend 1-color gray intensity processing to 3-colors [50], and are generally based on either histograms or clustering techniques. In histogram-based approaches, clusters are obtained by finding frequency peaks in the histogram. Unlike gray level histograms, color histograms have more than one dimension. Thus, the peak can be found either independently in each color channel or in the whole 3D histogram. In clustering-based techniques, pixels are grouped based on a distance metric and the color values of the points. The spread within a cluster is mainly determined by color variations due to shading and device noise. An example is the K-means algorithm, which iteratively computes each cluster's membership and mean color values until convergence. Weeks et al. showed that K-means provide efficient pixel classification based on color information [53].

### 3.1.1 K-means for color segmentation

In the K-means algorithm, an objective function, given below as  $J$ , is minimized.

$$J = \sum_{j=1}^K \sum_{i \in S_j} |x_i - \mu_j|^2 \dots\dots\dots (3.1)$$

In this case,  $J$  is the squared Euclidean distance of the  $n$  data points from their respective cluster centers, where  $x_i$  is a data point in cluster  $j$ , and  $\mu_j$  is the  $j^{\text{th}}$  cluster center. For color segmentation,  $x_i$  is a three dimensional vector of {red, blue, green}, but this can be generalized to any number of dimensions.

In general, the K-means algorithm suffers from the local minima problem and a lack of user control. In [22, 49, 54] the authors have shown that K-means clustering with user-provided seed points can result in better segmentation when compared to fully automated “unsupervised” K-means. Generally, clinical applications prefer accuracy over automation, if the segmentation results with user interaction can be produced in real time; they will be the top design choice. Under this guideline, i.e. to maximize the user feedback while performing the computation in real time, we have designed a color visualization and segmentation system.

This system maps the color information in 3D RGB space to 2D for faster interactive segmentation. Our 2D color space consists of the hue and chroma of the input colors. The user can mark zones on the color map rather than clicking on an individual seed point. This helps in better capturing the user’s perception, and results in better color segmentation.

### **3.1.2 Color segmentation using color maps**

The user’s perception is another very important aspect for any color segmentation technique. The color models used and the differentiation metric used to separate different classes should incorporate human perception models. A perceptual color segmentation algorithm [55] segments RGB color space into ten color categories using Munsell and CIELUV color models. This method has some user interaction since the users are asked to name representative colors in each color category, however, it is subjective due to different names given to the same color by different users. In [56] the authors find the dominant colors in an image by finding ridges in the color distribution and assigning a unique color at every ridge as a representative color of an interesting region. This approach has the limitation of labeling a region with a single color resulting in loss of color variation. Another method [57] is again related to skin color segmentation by using a 2D plane in RGB color space. Principal component analysis is applied based on the fact

that skin colors in the RGB space are approximately distributed in a linear fashion. The authors claim that the problem of color constancy [58] is relieved, but their results are for a specific particular dataset i.e. skin color images. Their algorithm may be applicable for some biological images in which the color distribution is linear in RGB space, but it may not work for other images in which the color distribution is not easily described by a plane in RGB space.

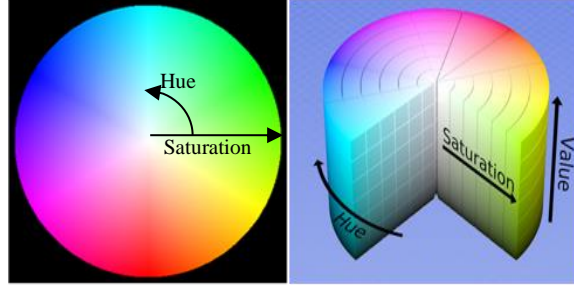
The user's perception of color can also be incorporated by introducing interactivity and user feedback. Interactive segmentation [59], in which the "user is in the segmentation loop," allows the addition and removal of regions of interest by the user. It does not, however, allow the user to select color components for segmentation. Interactive feature extraction has been applied to color breast cancer images [60]. In the training phase, that algorithm requires the user to select between 25 and 100 points for each color of interest to segment. This approach is highly subjective, especially when the colors for different image features vary only slightly from each other. In addition, this process is time consuming. The problem gets more complicated with color variation in biological images. A tool for interactive object segmentation in color images [61] has previously been developed. In that tool, seed point information is acquired from the user to perform watershed based segmentation. Another work [62] depends on user selected seed points for segmentation. The test image used in that work has clearly defined color edges and different objects in the image have homogenous colors, but most of the images used for biological purposes have diffused color boundaries and lack homogeneously colored objects. Another work [63] describes adaptive region growing color segmentation for text images. Authors have reported satisfactory results with text images except for images with very high color variance. This makes their approach unsuitable for biomedical image segmentation.

We designed our color segmentation tool to be interactive and to incorporate color constancy perceived by humans. Some other salient features of our segmentation scheme include:

- Visualization of 3D RGB colors in a 2D space that covers the pixel frequency, color composition and color class clustering (co-localization) of the images
- Selecting several color class samples in a time frame equivalent to single seed point selection
- Enhanced user control
- Compatibility with different types of cancer images stained with different colored dyes

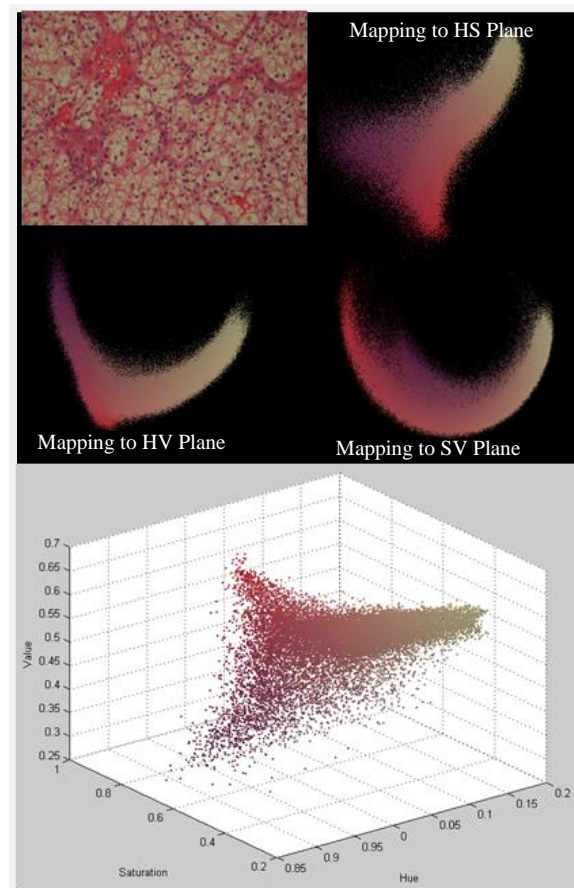
Segmentation in color space is challenging because it is not easy to find a similarity metric that translates the user's perception of 3D color space to a single dimension decision parameter. Even with high-speed computer graphics, it is difficult for users to select or encapsulate clusters directly in a 3D space, such as the RGB space. In this work, we reduce this complexity to provide an interactive environment in 2D. For data visualization we used Hue (H) and Saturation (S) plane (Figure 3.1) as mapping space where all the pixels in the test image are mapped close to other pixels which have similar HS values (co-localization). It may be noted that the color space is not being used to do the dimensionality reduction. It is only used to co-localize the pixels e.g. using HS plane, all pixels with same Hue and Saturation values will be mapped close to each other irrespective of intensity. The pixels displayed will still appear with their original intensity value i.e. the pixels in color map show all components (RGB).





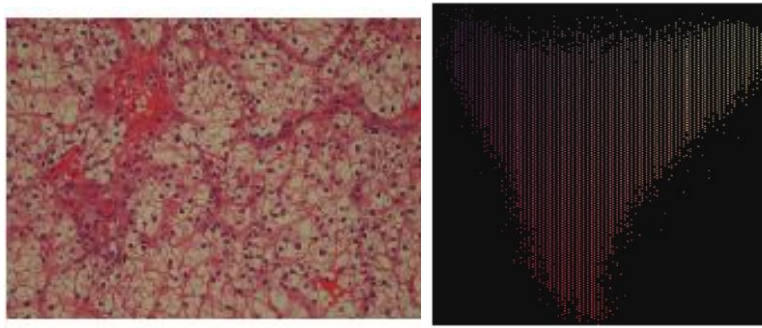
**Figure 3.1.** (L) HS plane with unit value (R) HSV color space[64]

We tested different visualization schemes using different images and selected HS plane for color maps as it showed better spread and easier interaction. One such example is shown in figure 3.2 comparing color maps based on HS, HV and SV planes (map generation explained later). It can also be seen that HS plane provides a better representation of test image pixels 3D scatter plot than HV and SV representation.



**Figure 3.2** Mapping of image pixels to different planes is shown using random spread visualization. Mapping to HS Plane presents better spread, easier interaction and better representation of 3D scatter plot of image pixels.

Direct mapping of the pixels to visualization space has certain limitations. Because most image formats support 24-bit color values, the probability of multiple pixels mapping to a single location is significant. In addition, if the visualization space image size is large in comparison to only 255 (8 bit) possible values along one axis, blank rows and columns remain, as shown in figure 3.3. We can also see that figure 3.3 does not give a true representation of how frequently a specific shade of color is represented in the input image.



**Figure. 3.3** (L) Test image (R) 2D Visualization space representing image. Both the problems of co-located pixels and the unmapped space can be seen.

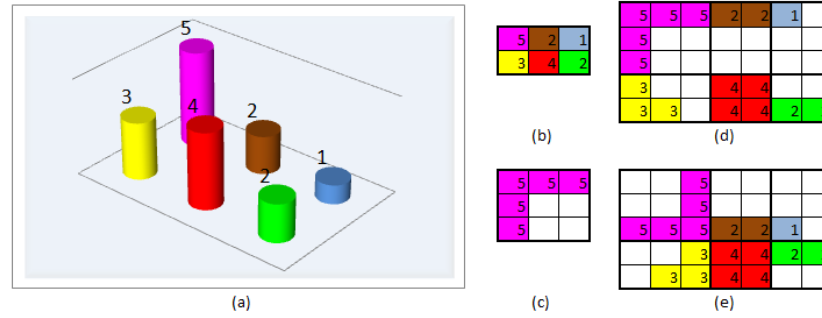
We tried different methods to generate color map by arranging pixels in order to present the user with a true representative of the image in real time for interactive visualization. In the following section, we will discuss different techniques, along with their advantages and disadvantages, to overcome these issues in order to present the user with a representative interactive visualization.

### Methodology Design

#### *A. Flattened Histogram Visualization*

2D visualization shown in Figure 3.3 loses the information of all the pixels which map to same location. If we retain this information we can generate a 3D histogram where the bin height represents number of pixels mapped to a particular location. To reduce dimensionality we can collapse these histogram bin towers by spreading them in

2D. This will result in a map where similar pixels are placed near each other and every pixel in the input image appears once in the map. We refer to this as a “Flattened Histogram.” This process involves creating sufficient space around each bin according to its height and then filling in the pixels data in these locations. There are lots of blank spaces created in this process which are subsequently removed by a radial sort scheme. Figure 3.4 demonstrates the flattened histogram process. Pixels of similar color in an image are mapped to a single location and pixel count is used to create a histogram as shown in Figure 3.4a. Figure 3.4b represents the same histogram in 2D visualization space and the number in each cell represents the number of pixels (histogram height) of that specific color in the input image. We compute a square space which can



**Figure. 3.4** (a) Histogram showing pixel count in input image binned according to the pixel color (b) 2D visualization of Histogram values (c) Space allocation for individual cell (d) Space allocation for complete map (e) Visualization map after radial sort

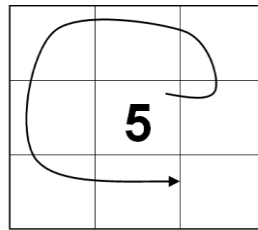
accommodate this number of pixels. For 5 elements we need a space  $\sqrt{5} = 2.23$ , rounding to next whole number we need a 3x3 space as shown in Figure 3.4c. By computing the row and column maximum values and computing the square space we generate the complete map as shown in Figure 3.4d. A radial sort is performed to move the pixel values towards the center and move the blank spaces outwards, towards the edges Figure 3.4e.

The resultant image shown in Figure 3.5b represents every pixel in the image (Figure 3.5a). Pixels are generally placed in the neighborhood of similar pixels.

However, two major problems were observed in this scheme. One is that during radial sort pixels may move significantly away from other pixels of their class. As shown in Figure 3.5b, some of yellow pixels in the flattened histogram have moved to the neighboring purple space. Secondly, the space allocation and radial sort are relatively expensive in terms of computational cost.



**Figure. 3.5** (a) Original Image (b) Flattened Histogram



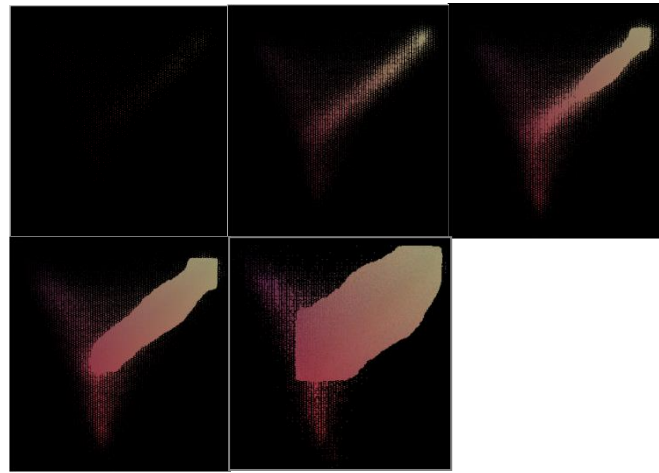
**Figure. 3.6** Spiral search

### *B. Spiral Spread Visualization*

Another approach “Spiral Spread” maps every pixel to an already occupied location, searches for an empty space in their vicinity and places it at the first available empty space. The empty space is found by searching in a spiral loop (Figure 3.6) increasing the search radius after completing each rotation. The process works fine in areas where the map is sparse but in places where spaces start getting filled, the search

radius starts increasing. The increased search radius is an undesirable condition which not only increases computation time but also maps pixels far away from their matching counterparts.

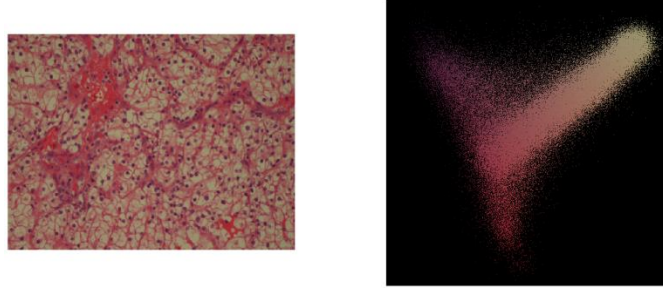
The spiral spread method represents pixel value frequency well, but suffers from a high computational cost as well as from similar pixels appearing far apart, especially in dense regions. Figure 3.7 shows different stages of the visualization map and its comparison with Figure 3.5 shows a better co-localization and better frequency representation of the input image pixels.



**Figure 3.7** Different stages of spiral spread scheme showing the evolution of map.

### *C. Random Spread Visualization*

In this scheme, rather than finding an empty location for every pixel, each pixel is mapped to a location with a random offset. The offset is small enough so that the pixels belonging to a similar class are mapped close to each other. Few pixels may overwrite others but the zones of frequently occurring pixels filling earlier than others can be seen as a measure of co-localization. Figure 3.8 below displays the original image and the spread of pixels in the corresponding visualization scheme.



**Figure 3.8** (a) Original Image (b) Random spread visualization of the original image

This visualization scheme is limited in displaying the true frequency of pixels, but it is computationally much faster than the two approaches already discussed. Another advantage is that it does not allow colors of one class to move far from their own class as they are constrained by the random distance limit. Area, along with the pixel density (compactness of color map), represents how frequently pixels of the similar color occur in the input image. For example, the scarcity of red and purple color in the input image, as compared to pink and white color, is demonstrated by the thick localization of pink and yellowish color in Figure 3.9(b).

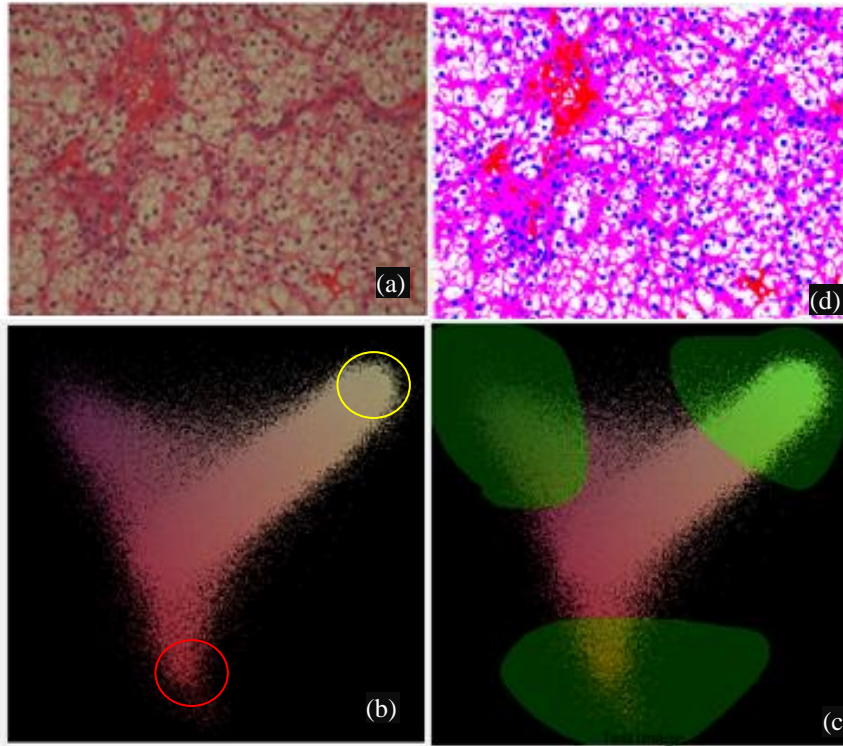
The performance statistics of these techniques are shown in table 3.1. Based on the performance and advantages offered by the ‘Random Spread Visualization’ over ‘Spiral Spread Visualization’ and ‘Flattened Histogram Visualization’, we selected ‘Random Spread Visualization’ as our default map for further processing.

Table 3.1 Performance comparison of different visualization techniques

Visualization Technique	Image size	Processing time (secs)
Direct mapping	800x600	0.279
Flatten histogram	800x600	169.654
Spiral spread	800x600	416.062
Random spread	800x600	0.493



The users are presented with this 2D color space visualization of (Figure 3.9b) to approximately mark the areas (Figure 3.9c), using the mouse, where pixels should be grouped together as a single color class. This zone marking is then used to segment the color image. The segmented image (Figure 3.9d) is completely dependent on the color ranges specified. Closed curves of any shape and size are acceptable and any number of color ranges can be specified, thus providing user flexibility.



**Figure 3.9** (a) Test image (b) Random spread visualization showing dense regions corresponding to frequently occurring image pixels(c) User marked zones based on different color classes in test image. (d) Pseudo color segmentation results.

To compare the segmentation results based on our color map we used level sets algorithm and few semi-supervised variants of K-means algorithm. We provide a brief description of these methods before presenting the segmentation results.

#### *Seeded K-means (SK-means).*

SK-means is our semi supervised variant for standard K-means[65]. Standard K-means is an unsupervised algorithm. Even initializing the algorithm with seeds points may help early convergence but resulting final classification may not be anywhere close

to starting seed points or the user's perception. SK-means utilizes user provided seed points for mean initialization as well as anchor point to restrict excessive mean movement to incorporate user's perception. During the iterative processes of mean computation, association of each point is determined based on its distance from new mean as well as starting seed point.

In SK-means, we split  $N$  data points  $x^{(n)}$ , in an  $I$  dimensional space, into  $K$  clusters with means  $m^{(k)}$  and initial seed points  $s^{(k)}$ . Each vector  $x$  has  $I$  components  $x_i$  and its distance from  $k^{\text{th}}$  mean and  $k^{\text{th}}$  seed point is given by

$$d_k^{(n)} = \sqrt{\sum_i (x_i^{(n)} - m_i^{(k)})^2} + \sqrt{\sum_i (x_i^{(n)} - s_i^{(k)})^2}$$

At start of SK-means, all means  $m^{(k)}$  are assigned seed point values  $s^{(k)}$  i.e.  $m^{(k)} = s^{(k)}$ . In the assignment step, each data point  $n$  is assigned to the nearest cluster based on minimization  $d_k^{(n)}$  which includes both the distance from the current cluster mean and the cluster seed point. The new assignment of clusters for all points is given by

$$\hat{C}(n) = \underset{k}{\operatorname{argmin}} \{d_k^{(n)}\}$$

All the data points  $x_{(n)}$  assigned to a cluster  $C$  form part of its responsibility and is given by the responsibility indicator  $r$

$$r_k^{(n)} = \begin{cases} 1 & \text{if } \hat{C}(n) = k \\ 0 & \text{if } \hat{C}(n) \neq k \end{cases}$$

In update step the updated means are computed as

$$m^{(k)} = \frac{\sum_n r_k^{(n)} x^{(n)}}{\sum_n r_k^{(n)}}$$

Repeat the assignment step and update step until the change in mean assignment is below a threshold.

### *Multi-seeded K-means (MSK-means).*

MSK-means is similar to SK-means except that each class is subdivided into multiple subclasses. This helps in better user feedback by providing more number of seed points (one per subclass). Intra-class variations can be captured by corresponding variation in



the selected seed points. After the classification process the subclasses are merged back to generate main class labels.

### *Level Sets Segmentation*

Level set method,[66] is an edge based technique to perform image segmentation. The boundary of the segmented object is defined as the zero level set of function  $\Phi(x, y)$ , i.e. it is implicitly defined as the solution of the equation  $\Phi(x, y) = 0$ . The boundary (and consequently  $\Phi(x, y)$ ) is initialized using multiple seeds placed at regular interval in a grid and then evolves until it conforms to the image. In order for the boundary to evolve,  $\Phi(x, y)$  has to evolve. To achieve this,  $\Phi(x, y)$  is added a time dimension, i.e. it becomes a function of three variables  $\Phi(x, y, t)$ . The equation (the level set equation) governing the change of  $\Phi(x, y, t)$  is

$$\frac{\partial \Phi}{\partial t} + F|\nabla \Phi| = 0$$

where  $F$  is the speed of the boundary in the direction normal to the boundary and

$$|\nabla \Phi| = \sqrt{\left(\frac{\partial \Phi}{\partial x}\right)^2 + \left(\frac{\partial \Phi}{\partial y}\right)^2}$$

Our implementation of level-sets uses intensity  $I(x,y)$  threshold and curvature  $k(x,y)$  to compute the speed function.

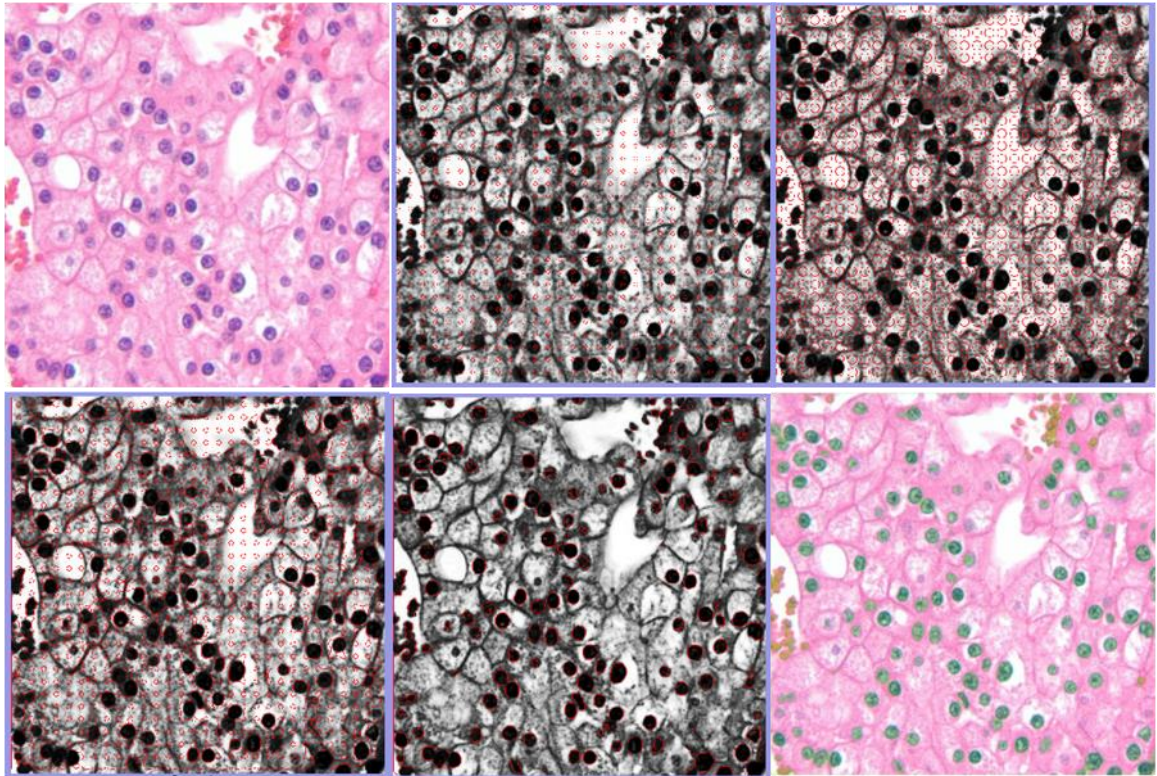
$$F(x, y) = f_1(x, y) - f_2(x, y)$$

$$f_1(x, y) = \begin{cases} 1, & \text{if } I(x, y) > \text{threshold} \\ 0, & \text{otherwise} \end{cases}$$

$$f_2(x, y) = \begin{cases} 0, & \text{if } k(x, y) < 0 \\ \frac{k(x, y)}{k_{max}}, & \text{if } 0 < k(x, y) < k_{max} \\ 1, & \text{if } k(x, y) > k_{max} \end{cases}$$

Where  $I(x,y)$  is the intensity and  $k(x,y)$  is curvature at a point while maximal curvature parameter is denoted as  $k_{max}$

An example of levelset segmentation is shown in figure 3.10.



**Figure 3.10** Level sets based segmentation is shown. Contours initialized at nodes of a grid, evolved and merged to generate segmentation mask.

### Weighted SK-means (WSK-means)

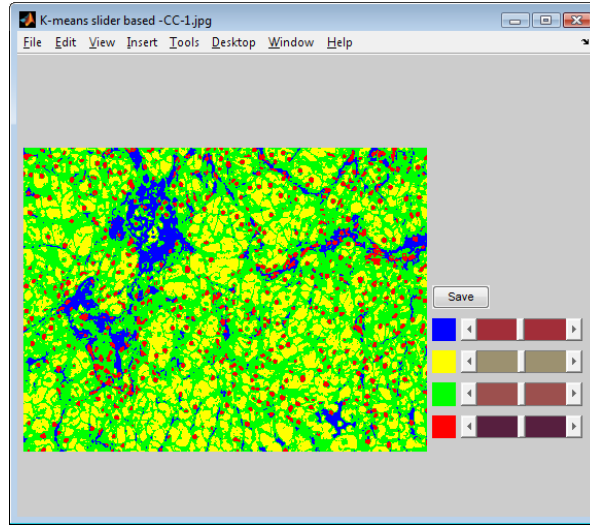
This variant of SK-means allows user to change the classification results with the help of sliders bars in a GUI. The original equation of SK-means is modified as

$$d_k^{(n)} = W_k \left( \sqrt{\sum_i (x_i^{(n)} - m_i^{(k)})^2} + \sqrt{\sum_i (x_i^{(n)} - s_i^{(k)})^2} \right)$$

where  $W_k$  is the weight parameter adjusted by the slider control.

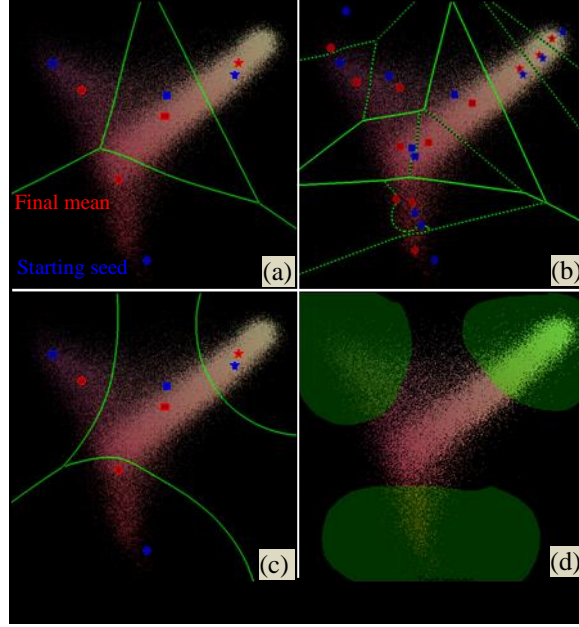
The user keeps manipulating the sliders controls visualizing the segmentation results in real time until desired results are achieved.

An example of WSK-means segmentation is shown in Figure 3.11.



**Figure 3.11.** GUI for WSK-means showing pseudo color representation of test image (figure 3.9). Sliders are used to change weights.

A better comparison of the aforementioned algorithms can be done using the visualization in figure 3.12. It may be noted that in actual practice all algorithms except color map utilize 3D space for classification.



**Figure 3.12** (a) Class partitions shown for SK-means segmentation. (b) MSK-means multiple sub-classes can be combined to cater for intra-class variations. (c) WSK-means can change classification result by increasing and decreasing class zones based on slider weights. (d) Interactive Color map – user can select zones of any shape and size.

### Testing and Comparisons

To best of our knowledge, no histological image data is available where each pixel is labeled that can be used as a ground truth for testing different algorithms. In the absence of such a dataset, we tested our algorithm for H&E, IHC stained biopsy tissue images of renal cell carcinoma and head & neck cancer. We devised a strategy to compare the accuracies of afore mentioned algorithms. We consider WSK-means algorithm results as ground truth since it provides interface to the user to update classification results until the best perceived segmentation is achieved. To support our decision, we used the level sets algorithm as the second reference. Our level sets implementation, based on the gray level intensity cost function, is suited best to segment the dark stained nuclei in the image. Comparing only this class we were able to show that accuracies above 95% can be obtained which are close to optimum as the two algorithms inherently differ from each other.

Having verified our ground truth results, we used them to conduct a human interaction study to compare the performance and accuracy of other algorithms under consideration. Our study involved eight users with varying background from new to frequent users of histology images. We used six H&E stained images of Renal Cell Carcinoma (RCC) images which were divided in two groups of 512x512 and 1024x1024 pixels to compare the variation in computational cost with image size. The users were presented with the previously described 2D color space visualization of (Figure 3.9b) and were prompted to approximately mark the areas where pixels can be grouped together as a single color class. For SK-means, the users mark seed points on the test image which best represents each class while in MSK-means user select multiple seeds per class to capture intra-class variation. The user interaction time including computational time of all the algorithms was recorded and all the algorithms were compared for performance and accuracy.

## Results

The performance results (Table-1 and Figure 3.13) show that random spread color map is faster than MSK-means and the average time taken by the user is less than 30 seconds per image, which is reasonable for most applications. Color map performance advantage, being a non-iterative algorithm, becomes apparent when the image size becomes large. The computational time for both SK-means and MSK-means time nearly doubled for larger images and only increased fractionally in case of color maps segmentation.

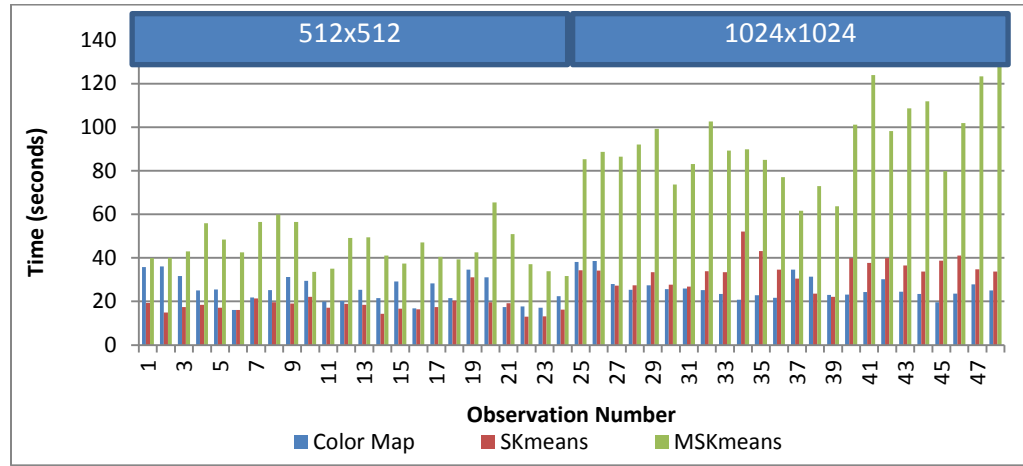
To compare accuracy of segmented images either region-based or pixel-based metrics are used [67]. In our image data we are interested in texture information inside the objects (nuclei) so we compared pixel labels from each segmented image. We used class accuracy, commission error and omission error as our evaluation metrics. The segmentation accuracy results are shown in table-2. As mentioned earlier WSK-means was used as ground truth validated by best level sets segmentation results selected after processing at various intensity thresholds. It can be seen that MSK-means performs better than SK-means as expected and color maps showed even better accuracies than MSK-means. Further analysis showed that accuracies are also dependent on type of the image as well as user's training and practice.

Table-1. User interaction time comparison

	Time (secs) $\mu \pm \sigma$			
	Level sets (single class)	Color maps	SK-means	MSK-means
Images 512x512	150.78 $\pm$ 46.01	25.04 $\pm$ 6.26	18.23 $\pm$ 3.60	44.87 $\pm$ 9.08
Images 1024x1024	1593.82 $\pm$ 39.77	26.48 $\pm$ 5.09	34.19 $\pm$ 6.79	91.28 $\pm$ 16.60

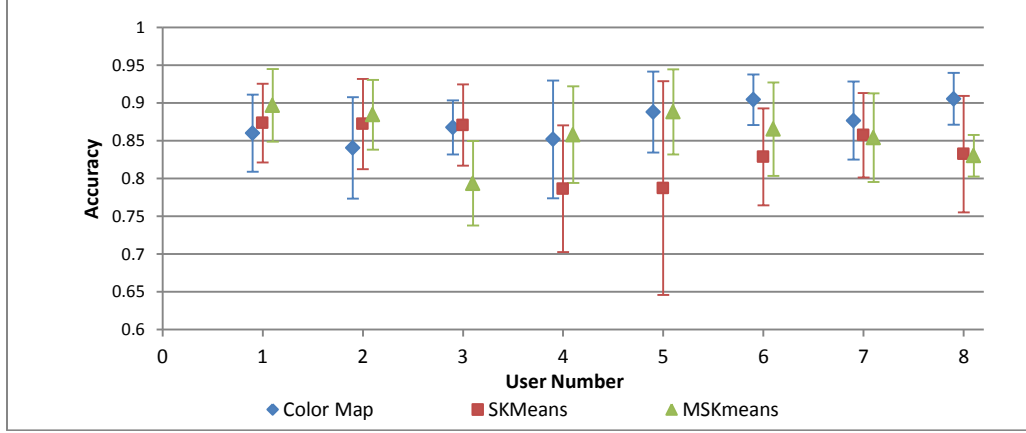
Table-2. Segmentation accuracy comparison

	Class Accuracy ( $\mu \pm \sigma$ )	Commissions ( $\mu \pm \sigma$ )	Omissions ( $\mu \pm \sigma$ )
<b>WSK-means</b>	1.0 $\pm$ 0.0	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0
<b>Color Maps</b>	0.874 $\pm$ 0.06	0.126 $\pm$ 0.06	0.126 $\pm$ 0.057
<b>SK-means</b>	0.839 $\pm$ 0.09	0.161 $\pm$ 0.09	0.161 $\pm$ 0.09
<b>MSK-means</b>	0.859 $\pm$ 0.09	0.141 $\pm$ 0.09	0.141 $\pm$ 0.09
<b>Level sets</b>	0.977 $\pm$ 0.02	0.083 $\pm$ 0.04	0.023 $\pm$ 0.02

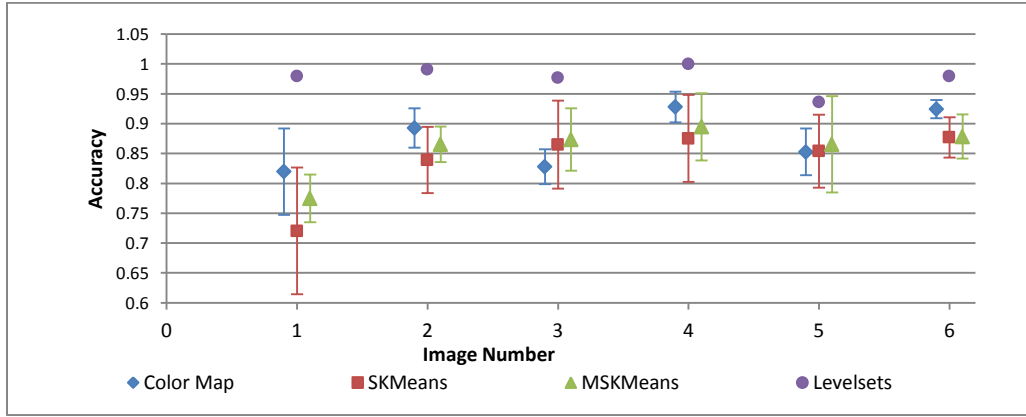


**Figure 3.13.** User interaction time (including computation time) comparison for two subgroups. Color map segmentation performs better for larger images.

Figure-3.14 shows results based on individual user. Users 6-8 were provided with a training session before performing actual segmentation. Users 1-5 were only given demonstration on single image before they actually performed segmentation. It can be seen that trained user can use the tool to get better accuracies especially in case of color maps based segmentation.



**Figure 3.14.** Accuracy comparison for individual users



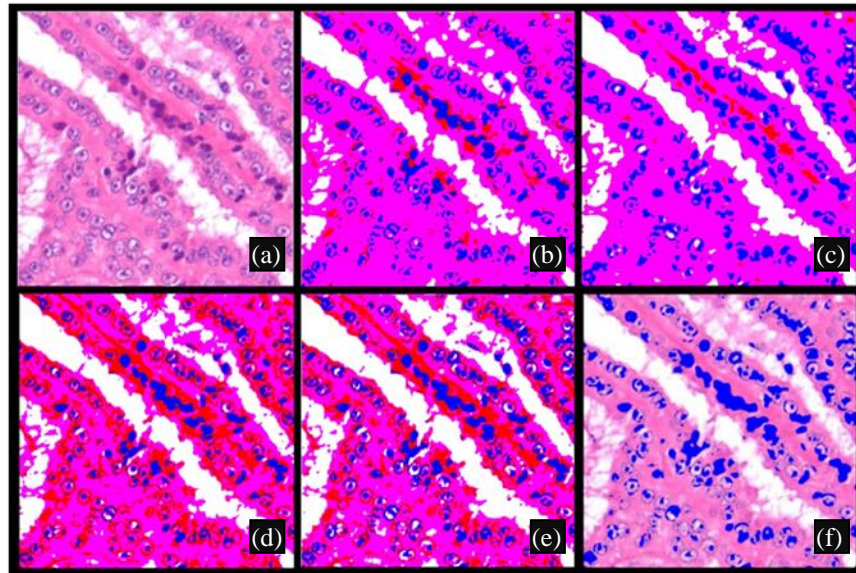
**Figure 3.15.** Accuracy comparison based on individual images in data set

Variation in accuracies can also be seen based on the type of image as shown in Figure 3.15. One of the images (Image-1 figure 3.15 and figure 3.16) has very small representation of one class (red). The control available in color maps helps in achieving better segmentation than other algorithms as highlighted in figure 3.16

We have tested our algorithm for H&E, IHC stained biopsy tissue images but our segmentation approach is applicable to other types of staining. . We were able to perform segmentation with 2% better accuracy in comparison to multi seeded k-means and in faster average time of 26 seconds on 1024x1024 pixel images in compared to 95 seconds in case of multi seeded k-means. In addition, our methodology presents important image statistics in a user friendly way. For example, the color spread provides three useful features of an image: its component colors, amount of each color and most importantly



co-localization or neighborhood information about color. Since we are presenting the 3D HSV color model with only two dimensions (Figure 3.1), we necessarily discard intensity in the model resulting in the co-localization of black, gray and white at the center of the visualization. Even with this limitation, we have shown that proper visualization of the component colors along with user interactivity results in better segmentation results. Our segmentation scheme is particularly useful for biological images since conventional schemes do not cater to the gradual color variations in these images. Moreover, the algorithm is iteration-free and executes faster than conventional K-means and its variants. Our segmentation tool is semi-automatic to capture miniscule color details. With our previous experience in the field, we see prospects of this work towards use of quantitative analysis and study of different types of histological images.



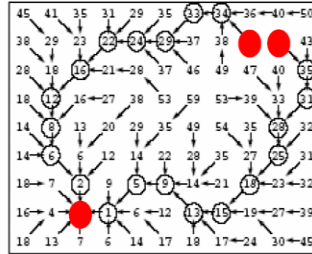
**Figure 3.15.** (a) Test image (b) WSK-means segmentation used as ground truth (c) Color map based segmentation (d) SK-means segmentation (e) MSK-means segmentation (f) Level sets based segmentation for single class overlaid on top of test image



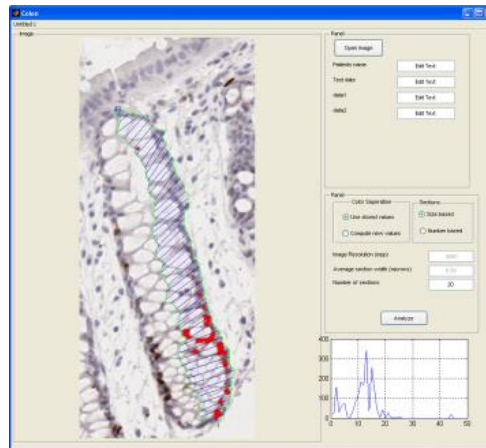
### 3.2 Region of Interest (ROI) segmentation

A ROI [68] is a portion of an image that you want to perform some specific operation on. You define an ROI by creating a binary mask, which is a binary image that is the same size as the image you want to process with pixels that define the ROI set to 1 and all other pixels set to 0. You can define more than one ROI in an image. The regions can be geographic in nature, such as polygons that encompass contiguous pixels but most of the times these are based on free hand user markings. In context of quantitative pathology, at times it becomes very critical and important to mark the regions very precisely and consistently for repeated experiments. This necessitates some intelligent algorithms to be used for marking ROI. In addition, if it facilitates the user in terms of speed and ease of use. We demonstrated one such technique for segmenting Colon crypts [42] or more precisely hemi-crypt (one half of the symmetric colon crypt). This is accomplished by using our semi-automatic segmentation tool based on the intelligent scissor (IS) algorithm [69, 70]. IS is a novel approach in object segmentation. Rather than optimizing a user-initialized approximate contour, IS allows the user to interactively select a boundary from a collection of optimal solutions. IS interactively computes the optimal path from a user selected “seed” point to all other points in the image. The optimal path from each pixel is determined at interactive speed by computing an optimal spanning tree of the image using an efficient implementation of Dijkstra’s graph searching algorithm. This search algorithm treats the image as a weighted graph (Figure 3.16). Each pixel represents a node with directed and weighted edges that connect with its eight adjacent neighbors. As the cursor moves, the optimal path from the pointer position to the seed point is displayed. This path allows the user to select an optimal contour segment that visually corresponds to a portion of the desired object boundary. As the mouse pointer comes in proximity to an object edge, a live wire boundary snaps to and

wraps around the object of interest. An example of this algorithm is shown in Figure 3.17.



**Figure 3.16** An example of cumulative cost and path matrix used for optimal path computation



**Figure 3.17** An application snapshot showing Hemi-crypt segmentation using live-wire segmentation

### 3.3 Nuclear cluster segmentation

In pathological conditions, complex cell clusters are a prominent feature in tissue samples. Segmentation of these clusters is a major challenge for development of an accurate cell counting methodology. We address the issue of cluster segmentation by following a three step process. The first step involves pre-processing required to obtain the appropriate nuclei cluster boundary image from the RGB tissue samples. The second step involves concavity detection at the edge of a cluster to find the points of overlap

between two nuclei. The third step involves segmentation at these concavities by using an ellipse-fitting technique. Once the clusters are segmented, individual nuclei are counted to give the cell count. The method was tested on four different types of cancerous tissue samples and shows promising results with a low percentage error, high true positive rate and low false discovery rate.

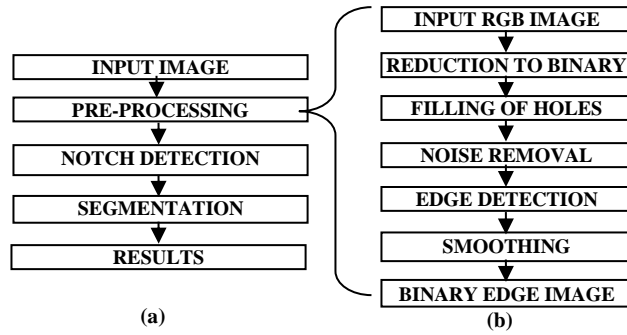
Pathologists often depend on parameters such as the number, shape and size of cells in a tissue sample to make important diagnostic decisions. In healthy conditions, nuclei in cells are mostly distinct and parameters can be determined by direct image segmentation methods such as region-based methods, histogram-based methods and edge detection based methods. However, in pathological conditions, individual cells come close together and nuclei form dense clusters. Figure 3.19(a) shows dark elliptical nuclei touching and overlapping in a 2-D tissue sample. Therefore, accuracy of cell-counting, cell shape and size determination depends on the segmentation of these dense clusters.

Previous work addresses segmentation of simple-clusters and touching cells by extending and improving image-segmentation methods [71, 72]. Few authors have developed algorithms that address cluster segmentation specifically [73-76]. All these methods addressing cluster segmentation either could segment only simple clusters [71, 72] or give good results only for circular cells [71, 72, 76] or resulting cell shape is not a good model for the original cell shape [73-75] or have very complex algorithm [74-76]. However, this paper presents an edge-based image segmentation method, shown in the flow-diagram (Figure 3.18(a)), that is simple to implement and can segment complex clusters with reasonable accuracy. The method involves detection of concavities on cell cluster edges and segmentation at these concavities by ellipse fitting. The elliptical model used is a good approximation to the original cell shape. Recently, Wang and Song [77] and Bai et al. [78] introduced the concept of cluster segmentation using concavity detection. In this paper, we present a novel method for notch detection using cross-product (section 3.3.2), and a new technique for cluster segmentation using ellipse fitting

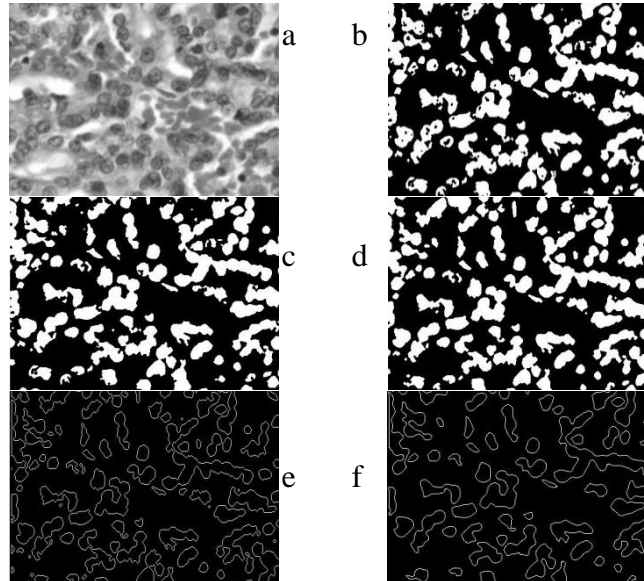
(section 3.3.3). Also, we perform quantitative analysis of segmentation result using standard statistic parameters. Using proposed methodology, pathologists will be in better position to take diagnostics decisions.

### 3.3.1 Preprocessing

We have implemented our method for different types of tissue samples including standard photo micrographs of H&E stained biopsy tissue sections of renal cell carcinoma (RCC) and IHC stained biopsy tissue sections of head and neck (H&N) cancer.



**Figure. 3.18** a) Overall flow-diagram for the method, b) Pre-processing steps



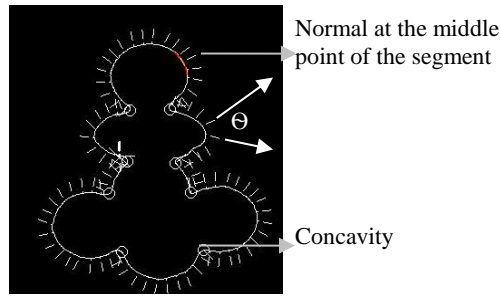
**Figure. 3.19** Pre-processing steps implemented for papillary tissue sample. (a) Input RGB image shown in gray scale, (b) binary mask of nuclei, (c), filled binary mask, (d) mask after noise removal (e) result after edge detection, (f) result after smoothing

Due to the nature of tissue images and the variability in the sample preparation, staining, and image acquisition process, it is imperative to pre-process these images in order to remove variations. The pre-processing steps have been depicted in the flow-diagram in Figure 3.19(b) and the corresponding images are shown in Figure 3.19. In stained RGB tissue images various entities in a tissue slice such as nuclei, glands, cytoplasm and red blood cells appear as different colors. The first pre-processing step involves the generation of a binary mask for cell nuclei from the RGB image using K-means clustering [53, 79], where seed points are selected by user interaction. The binary mask of a tissue sample often has clusters with holes as shown in Figure 3.19 (b). If these holes are not filled, they can be detected as false boundaries during the edge-detection process. Therefore, the next step involves filling in the holes using an algorithm based on morphological reconstruction [80] to obtain properly connected clusters as shown in Figure 3.19 (c). Very small objects in the binary mask are generally due to noise and due to misclassification during the k-means clustering. As such, the next step involves noise removal using the size threshold. Images in Figure 3.19(c) and Figure 3.19(d) show the mask before and after noise removal. Based on connected component analysis each object in the image is processed as an individual cluster. The boundary of each cluster is then detected based on a neighborhood of 8 pixels. Figure 3.19(e) shows result after edge detection. The resulting sequence of pixels that form the boundary of the cluster is then processed using smoothing techniques for better notch detection. Noise or jaggedness present on the edges of the clusters may lead to false concavity detection on the edge and consequently may be treated as a notch for segmentation. Therefore, it is necessary to make the boundary smooth and preserve true concavities. Our algorithm performs simple smoothing using a moving average low-pass filter. The Resulting image after smoothing is shown in Figure 3.19(f).

### 3.3.2 Concavity or notch detection

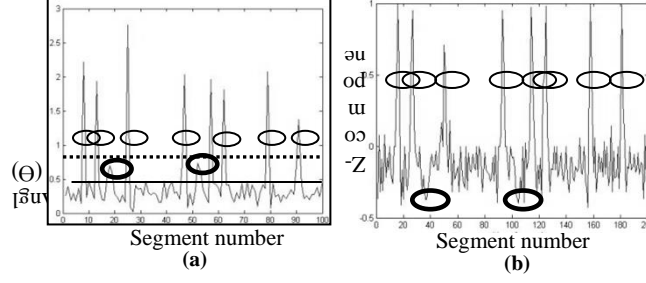
After preprocessing, the next step involves detection of concavities or notches. A concavity is the point on the cluster edge where two individual cells overlap. Therefore, the concavities can be found using angle ( $\theta$ ) between adjacent normals on the edge of the cluster as suggested by previous work [78]. In this method, we divide the edge of the cluster into fixed length segments, and plot a normal at the middle point of every segment as shown in Figure 3.20. If  $\Phi_i$  gives the slope of normal at middle point of segment  $i$  with respect to positive x-axis, then  $\theta$  for any segment  $i$  is given by:

$$\theta = \begin{cases} |\Phi_i - \Phi_{i-1}|, & \text{if } |\Phi_i - \Phi_{i-1}| < \pi \\ \pi - |\Phi_i - \Phi_{i-1}|, & \text{else} \end{cases}$$



**Figure 3.20** A synthetic cluster illustrating the method of calculating  $\Theta$ , angle between adjacent normals.

Edges have a depression around a concavity and a sudden change in surface orientation. Hence,  $\theta$  has maxima at concavities as shown in Figure 3.21(a). As illustrated in the graph using a high threshold (dotted line) only major concavities are discovered. These concavities are sharp concavities and in order to discover relatively smooth concavities the threshold needs to be decreased (solid line) and with this decrease some false detections start appearing at points with sufficient angle change, such as the ones at the edge of the individual elliptical cell with high eccentricity.



**Figure 3.21** a) Graph illustrating angles between adjacent normals for different segment number along the edge of cluster in figure 3. Dotted line and complete line represents high and low threshold respectively. b) Relation between z-component and segment number for the same cluster. Thin circles mark true concavities and thick circles mark the possible false concavities

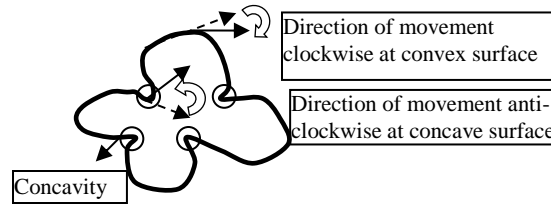
To avoid these false detections with a decrease in threshold, we exploit the fact that the desired concavities are located at the edge where the surface is concave (when viewed from inside the cell). As such, any detection at locations where the surface is convex can be rejected. The process involves splitting the cluster edge into segments of equal length. Vectors are generated for tangents at every segment. The cross product of each pair of adjacent tangential vectors is calculated while moving in clockwise direction along the cluster edge. The magnitude of the cross product depends on the angle between the vectors and its sign depends on the direction in which the first vector moves towards the second vector.

$$\mathbf{a} \times \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \sin \theta \hat{\mathbf{n}} \quad (2)$$

Where,  $\mathbf{a}$  &  $\mathbf{b}$  are first (dotted line) and second (solid line) tangential vectors;  $\theta$  is the angle between the vectors, and  $\mathbf{n}$  is the unit vector perpendicular to  $\mathbf{a}$  and  $\mathbf{b}$  in the direction given by right hand rule.

If cluster is in X-Y plane, Z component of the cross product represents magnitude and direction of cross product. As shown in Figure 3.22 at convex surface, the direction of movement is clockwise and  $\mathbf{n}$  is negative z-direction. While at concave surface, the direction of movement is anti-clockwise and consequently  $\mathbf{n}$  is positive z-direction. High positive z component represents notch while negative value represent convex surface. Comparing Figure 4.6 (a) and Figure 4.6 (b), it can be observed that false detections in

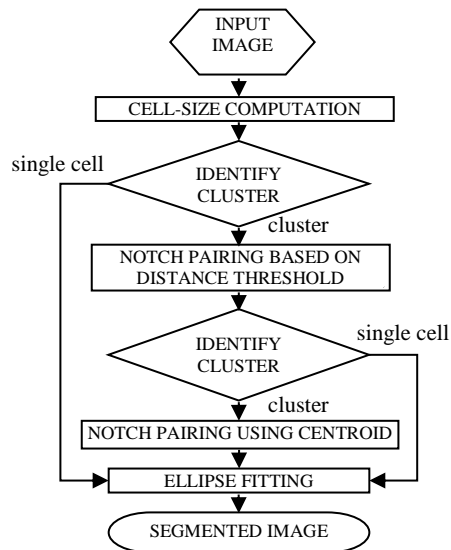
Figure 3.21 (b) are no longer affecting the concavity detection which can now be done at relatively lower threshold.



**Figure3.22** Figure depicts cross product resultant direction in case of convex and concave contour locations; concavities are marked

### 3.3.3 Cluster Segmentation

After detecting the concavities, we segment the cluster into cells. The segmentation algorithm assumes that cells have approximately elliptical shape with different eccentricities, suggested by a various authors. Figure 3.23 shows the flow chart for the segmentation of clusters.



**Figure 3.23** Flow-chart for Segmentation



### **3.3.4 Cell-size computation**

First, the average cell size is computed based on user interaction. The user selects a few samples of single cells. We then compute the average cell size from samples and use it to set a threshold for identifying cell clusters.

### **3.3.5 Cluster identification**

This step differentiates between cluster, single cell and noise. Depending on the standard cell size, two thresholds are set. The first threshold decides if any region is large enough to be treated as a cluster, and the second threshold decides if it is small enough to be treated as noise. The detected clusters undergo further segmentation while single cells are passed to ellipse fitting algorithm. Due to these thresholds the methodology is robust in segmenting cells in a tissue image with cell size within a range from the average cell-size.

### **3.3.6 Notch pairing based on distance threshold**

In this step, we compare the distance between all the notches. Any two notches that are closer than a particular threshold, which depends on average cell size, form a pair and the cluster is split at these two notches; preference is given to the notches that are closer. Cluster splitting based on distance splitting is continued iteratively until the point is reached where no further segmentation is possible. During this process each cluster segments into either sub-clusters or sub-clusters and individual cells. The sub-clusters generated during this step generally have circular shape and cannot be further split using the threshold criteria. Therefore, these sub-clusters are passed to next step for segmentation based on centroid.

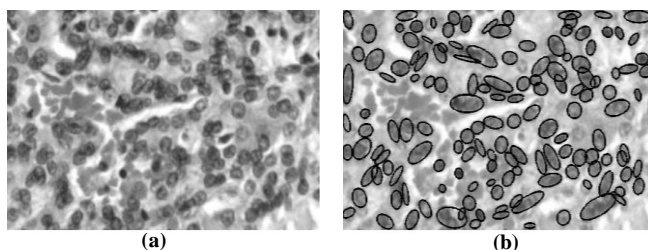
### **3.3.7 Notch pairing using centroid**

After distance-based segmentation, if there are any clusters left, they are segmented by using centroid connection. In this step, starting with the notch with highest

z-component (section 3.3.2), notches are connected through the centroid to split the cluster into cells. Any segmentation step is possible only if it results in regions with size larger than a minimum cell size threshold.

### 3.3.8 Ellipse fitting

We picked up the ellipse fitting method proposed by Fitzgibbon et al. [81] which is reported to have better accuracy than other standard methods. The edge pixels of cells obtained in the segmented mask after step 3.3.1 are compared with the edges of the cluster as shown in Figure 3.19(f). Common edges pixels are then used as data for ellipse fitting algorithm. The final result of the segmentation of the input image in Figure 3.24(a) is shown in Figure 3.24(b); black lines mark the cell boundaries.



**Figure 3.24** a) Input papillary tissue, b) result image after segmentation of image, green line mark the cell boundaries

### 3.3.9 Results

In order to test the robustness of our algorithm, we selected H&E stained tissue images from three subtypes of renal cell carcinoma (RCC) – papillary (PA), chromophobe (CH) and clear cell (CC) and IHC stained head and neck (H&N) cancer tissue images, thereby introducing morphological structure variations (RCC subtypes) and stain color variations (H&E and IHC). Quantitative analysis of these four different types (PA, CH, CC, H&N) of tissue images, numbered 1-4 respectively, is shown in table 3.3. Estimated number (EN) was calculated by the algorithm and false positive (FP) and False negative (FN) were estimated by comparison with manual segmentation results. The results have been analyzed based on three standard statistics parameters – True

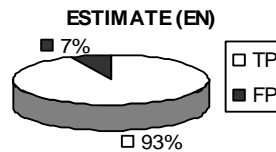
positive rate, false discovery rate and percentage error. High true positive rate, low false discovery rate and low error rate illustrate the usefulness of the method for cell-counting of various tissue samples. The method is simple to implement and can generate results in real-time, this is highly suitable for clinical applications. The method is semi-automatic and requires user interaction only in seed selection (section 3.3.1) and cell size calculation (section 3.3.4).

As compared to previous methods of segmentation using concavities [74, 77, 78], our method will generate better results due to higher accuracy in concavity detection. The method may generate errors at places when two cells overlap in such a fashion that concavities are very smooth or absent. Also the results are dependent on how good color segmentation is performed during the k-means clustering process for generation of the binary mask. The future work includes efforts to improve the color-segmentation and the cluster segmentation to further enhance the efficiency of the method. Also method is being tested for larger dataset of about 100 images to evaluate the robustness of the method.

Table 3.3 Quantitative analysis of nuclear segmentation for four different types of RCC

Image number	EN	FP	FN	TP	AP	TPR	FDR	E
1.	301	20	43	281	324	86.72	6.64	-7.09
2.	423	32	33	391	424	92.22	7.5	-0.23
3.	338	40	27	298	325	91.69	11.8	4.00
4.	480	60	22	420	442	95.02	12.5	8.59

Acronyms used in these tables are as follows: Estimated- Number of cells segmented by the method, EN; False positive- false detection, FP; False negative- missed detection, FN: True positive- correct detection, TP= EN – FP; Actual positive-number of cells calculated manually, AP= TP+FN; True positive rate- hit rate/sensitivity, TPR=TP/AP \*100; False discovery rate, FDR=FP/EN \*100; Percentage Error, E = (EN-AP)/AP \*100



### 3.4 Summary

In this chapter we reviewed some of the important segmentation techniques related to various segmentation domains. We also showed some new techniques and improvements over existing methods which can really be helpful for pathological image analysis tools. Color map based user interface, intelligent scissors based precise ROI segmentation and complex nuclei cluster segmentation are the techniques which proved to be really useful in achieving consistency and speed in the pathological image analysis.

## **CHAPTER - IV**

### **FEATURE EXTRACTION, SELECTION AND CLASSIFICATION**

The terms of feature extraction, selection and classification are widely used in the domain of pattern recognition. We will review their basic definitions and then see their application to our specific domain of pathological image analysis.

Feature extraction [82] is a special form of dimensionality reduction frequently used in pattern recognition. When the input data to an algorithm is too large to be processed then the input data will be transformed into a reduced representation set of features (also named features vector). Transforming the input data into the set of features is called feature extraction. If the features extracted are carefully chosen it is expected that the features set will extract the relevant information from the input data in order to perform the desired task using this reduced representation instead of the full size input. Best results are achieved when an expert constructs a set of application-dependent features. Nevertheless, if no such expert knowledge is available general dimensionality reduction techniques like Principal components analysis may help.

Feature selection [83] is the technique of selecting a subset of relevant features for building robust learning models. By removing most irrelevant and redundant features from the data, feature selection helps improve the performance of learning models by alleviating the effect of the curse of dimensionality, enhancing generalization capability, speeding up learning process and improving model interpretability. Feature selection also helps people to acquire better understanding about their data by telling them which are the important features and how they are related with each other. Feature selection algorithms typically fall into two categories: feature ranking and subset selection. Feature ranking ranks the features by a metric and eliminates all features that do not achieve an

adequate score. Subset selection searches the set of possible features for the optimal subset.

Classification [84] is the problem in statistics of identifying the sub-population to which new observations belong, where the identity of the sub-population is unknown, on the basis of a training set of data containing observations whose sub-population is known. Thus the requirement is that new individual items are placed into groups based on quantitative information on one or more measurements, traits or characteristics, etc. and based on the training set in which previously decided groupings are already established. The problem here may be contrasted with that for cluster analysis, where the problem is to analyze a single data-set and decide how and whether the observations in the data-set can be divided into groups. In certain terminology, particularly that of machine learning, the classification problem is known as supervised learning, while clustering is known as unsupervised learning.

In specific relation to pathological image analysis we will review the classification problem using different features sets for RCC applications [49, 54, 85] including knowledge based features, morphological, textural and wavelets based features and cellular features of elliptical models of segmented nuclei clusters.

#### **4.1 Knowledge based features**

Traditionally, most automated cancer diagnosis research has been on the problem of identification of cancerous and normal tissue images. Since pathologists use deviations in cellular structure as a means to make a diagnosis, many of the previous research efforts have used the statistical variation of various image properties to help make a diagnosis. The use of morphological features, for example, was reported by Jiang et al. in their study of breast cancer classification and by Roula et al. for the grading of prostate cancer [8, 9]. The diagnosis system developed by Diamond et al. used a combination of structural and

textural features to achieve an accuracy of 79.3% for the classification of Prostatic Neoplasia [10].

Esgiar et al. studied the classification of colonic mucosa using six different textural features and optical density and reported an overall accuracy of 90.2 % [11]. Their choice of features was motivated by the hypothesis that the human eye uses these features for texture discrimination. They reported an increase in accuracy of their system when fractal analysis was employed along with textural features and suggested the need for knowledge incorporation for further increase in accuracy [12]. Hamilton et al. used knowledge guided segmentation to calculate features like the co-occurrence matrix and optical density to study colorectal dysplasia achieving 83% correct classification [13].

Many of these studies have relied on blind self-training by selecting textural, morphological, topological or intensity based features or a certain combination of these based on the properties of the images under study. Although these various combinations of these features have been proven useful for cancer diagnosis of different types, their use will lead to better classification if employed in conjunction with prior structural properties of tissues under study.

We contend that by incorporating knowledge from an expert pathologist at every step of the system (Image processing, feature extraction, classification) the classification accuracy can be increased. Also by involving the user into the decision making process and allowing him to bias the system will lead towards making accurate prediction. Tissue images show both intra and inter class variation in terms of irregularity of cellular structure. Hence while CAD systems try to extract and quantify the inter class variations they should be adaptable enough to neglect the intra class variation. Having user interaction through every step of the process helps to encompass the vast non homogeneity that tissue image display and make the system more robust.

In this chapter we will present the system design, development and results of novel CAD based diagnosis system which allows an expert user to interact with the

system throughout the diagnosis process. Not only can the user bias and validate the results of feature extraction and quantification, he can select from a list of features he deems most important and appropriate for the classification. With such user interactivity and flexibility the same CAD tool can be used by pathologist from different cancer specialization to classify their images as long as the system has had sufficient training against images of that cancer subtype.

We chose RCC (section 1.3) as a case study for the development of this tool primarily because not much research has been done for the automated classification of renal tumor. Moreover this problem is much more complicated than normal/cancerous tissue classification as RCC has 4 common subtypes. Renal tumor subtypes exhibit several common morphological characteristics, making diagnosis difficult and subjective in many cases. Histopathologic classification is critical for the treatment of RCC as its histopathological subtypes are associated with distinct clinical behavior. So a diagnosis technique based on quantitative approach to renal tumor classification is critical and very much needed.

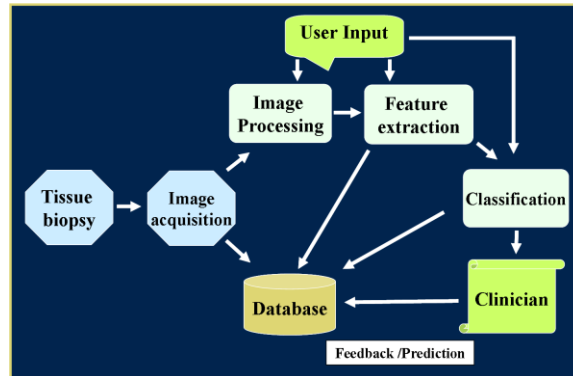
Expert knowledge of RCC features for feature extraction was incorporated into the system thereby increasing the classification accuracy. This was done by letting the pathologist select the features most relevant for the diagnosis of RCC. This was coupled with the prior knowledge about the presence and/or absence of specific histological features and structures (Red blood cells, blood vessels, lipid structures, papillary bodies) in various subtypes of RCC. The outline of the design of this novel diagnostic tool targeted for clinical practice and translational research is given in the next section.

#### **4.1.1 System Design**

The aim of the study is to take the past research advances in the field of quantitative molecular pathology one step further and develop this tool as a real clinical application. The whole design of the system has been built around that very goal. The



detailed outline of system proposed for system is given in Figure 4.1. It should be noted in particular that the user can interact and correct (if necessary) if necessary the results various image processing modules in addition to selecting various options. As mentioned earlier the whole point in providing this user interactivity is to keep the tool as general as possible and let the user incorporate his knowledge by correcting or biasing the results or by selecting between various options.



**Figure 4.1** Workflow of the proposed system showing data flowing between image acquisition, processing, feature extraction, classification, feedback and storage modules.

The images, once acquired, are processed to improve image quality and enhance the objects of interest. The regions of interests are then segmented out and passed onto the next stage for feature extraction and quantification. Features are extracted based on the expert knowledge built into the system about various expected RCC features. Essentially, the system tries to quantify features that describe the known difference among various classes. The images are then classified based on the extracted features and the results are provided to the pathologist for evaluation and feedback. The eventual goal of the study is to have a final system where the RCC images, the extracted features, the classification results along with the clinician's diagnosis are stored in a database. The database is used to train the classification system and image and feature annotation. This all ties in to the concept of having a practical and useful clinical tool for cancer diagnosis.

The following sections describes in detail the methods used for image processing, feature extraction and classification.

#### **4.1.2 Tissue Samples and Image Collection**

All tissues in this study were derived from renal tumors resected by total nephrectomy. Tumors were fixed, processed, sectioned and stained according to standard pathological procedures. Nephrectomy specimens were fixed for at least one hour in several volumes of 10% neutral buffered formalin, after which representative histologic samples (3-millimeter thickness) were obtained and fixed overnight in > 10 volumes of 10% neutral buffered formalin. Histologic samples were embedded in paraffin and microscopic sections (5-micrometer thickness) were prepared with a microtome and stained with hematoxylin & eosin. Representative photomicrographs of renal tumor sections were obtained at 200x total magnification and images of 600 x 800 pixels were extracted from the original 1200 x 1600 images for analysis.

#### **4.1.3 Image Processing**

Before extracting cellular/nuclear level information, image quality is improved by the image processing module. One of the main problems faced during segmentation of cellular features is the variation in staining. Both intra and inter image staining variation is observed in these biopsy images. In order to reduce this variation the images were first passed through a Gaussian smoothing filter. The size of the Gaussian filter used, is based on the average size of cellular features but the user has the ability to change the size of the filter if necessary. Although the Gaussian filter tended to blur the edges, it also smoothen the variation in staining. Segmentation results for filtered images were better than their unfiltered counterparts thereby validating the use of the filter. It is important to know that the image is smoothed only for the segmentation process. Once the objects within the image are recognized and tagged the pixels corresponding to those objects in

the original image can be used where necessary. This is particularly important for the case of textural feature extraction where a Gaussian smoothing would have altered the texture within the image.

Cellular/nuclear information present in the images is extracted in the next step. Both region-based and edge-based segmentation techniques were employed. Edge-based segmentation methods tended to suffer more at the hand of staining variation. The reason for this was a lack of significant intensity variation between the nucleus and its surrounding. Since the images were colored, color information in the images was used as a criterion for segmentation through K-means clustering algorithm and using user interactive tool discussed earlier in section 3.1.1

#### **4.1.4 Feature Extraction**

Segmented regions of interest, are now used to quantify cellular properties. Once quantified, it is these variations in cellular structure and distribution that are used for classification. Instead of calculating random features from the images or to find random hidden patterns, the system tries to model expert pathologist knowledge by extracting features corresponding to the morphological properties that were known to be different among various subclasses. These features were selected by the pathologist beforehand. In essence prior knowledge about expected image properties guides the selection of features to be extracted and the regions over which those features are calculated. An example of this is the papillary or finger-like feature present extensively in the papillary RCC. Fractal dimension is used to model this finger like structure. Fractal dimension has been used extensively in research to quantify the self-similarity of the images and is usually calculated over the whole image with or without various thresholding steps. By incorporating user knowledge the fractal dimension in this case was calculated over these fingers like structures only. Since this feature is inherent to papillary RCC, its values serves as a good descriptor for PAP.

The knowledge that different subtypes of RCC display variation in nuclei density and nuclei shape is captured by calculating various morphological features. PAP RCC tends to show higher nuclei to cytoplasm density while nuclei in CHR RCC images are known to be more eccentric than others and have a halo around them. Features like area, eccentricity, compactness ratio as well as ratio of the area of the nuclei to cytoplasm area are calculated to quantify these differences.

In addition to accounting for modeling the finger like structures in PAP, fractal dimension was calculated to quantify the self-similarity among features. This was based on the fact that cellular features, no matter how irregular in shape, display a level of self-similarity. Schepers et al. showed that among the various algorithms proposed in the past for the calculation of fractal dimensions, spectral analysis provides results with highest fidelity. Fractal dimension was calculated using both the spectral analysis and a box-counting method. Although the latter is a less complex algorithm to implement, the values from spectral analysis method were used because of their accuracy.

Table 4.1 Feature extracted for Papillary (PAP), Clear Cell (CC), Chromophobe (CHR) and Renal Oncocytoma (ONC) with mean and standard deviations

	PAP	CC	CHR	ONC
Correlation	8.23 $\pm$ 0.906	10.85 $\pm$ 0.971	2.505 $\pm$ 0.320	5.885 $\pm$ 1.01
Contrast	0.648 $\pm$ 0.077	0.43 $\pm$ 0.038	0.532 $\pm$ 0.031	0.454 $\pm$ 0.050
Energy	0.017 $\pm$ 0.005	0.018 $\pm$ 0.002	0.076 $\pm$ 0.016	0.029 $\pm$ 0.006
Homogeneity	0.504 $\pm$ 0.028	0.474 $\pm$ 0.011	0.663 $\pm$ 0.018	0.531 $\pm$ 0.025
Entropy	7.766 $\pm$ 0.047	7.446 $\pm$ 0.033	6.716 $\pm$ 0.071	7.288 $\pm$ 0.080
Fractal Dimension	1.787 $\pm$ 0.08	1.934 $\pm$ 0.007	1.842 $\pm$ 0.031	1.93 $\pm$ 0.022
Ratio of Area	0.796 $\pm$ 0.147	0.260 $\pm$ 0.035	0.105 $\pm$ 0.036	0.212 $\pm$ 0.073
Eccentricity	0.812 $\pm$ 0.025	0.827 $\pm$ 0.021	0.792 $\pm$ 0.052	0.766 $\pm$ 0.033

The knowledge about presence of RBC's and blood vessels in some subtypes (CHR), the presence of fibrous and vascular core (PAP), growth of nuclei in nest like structures (ONC) was quantified by calculating textural features like contrast, homogeneity, correlation and energy. RBCs are not stained by H&E staining and are

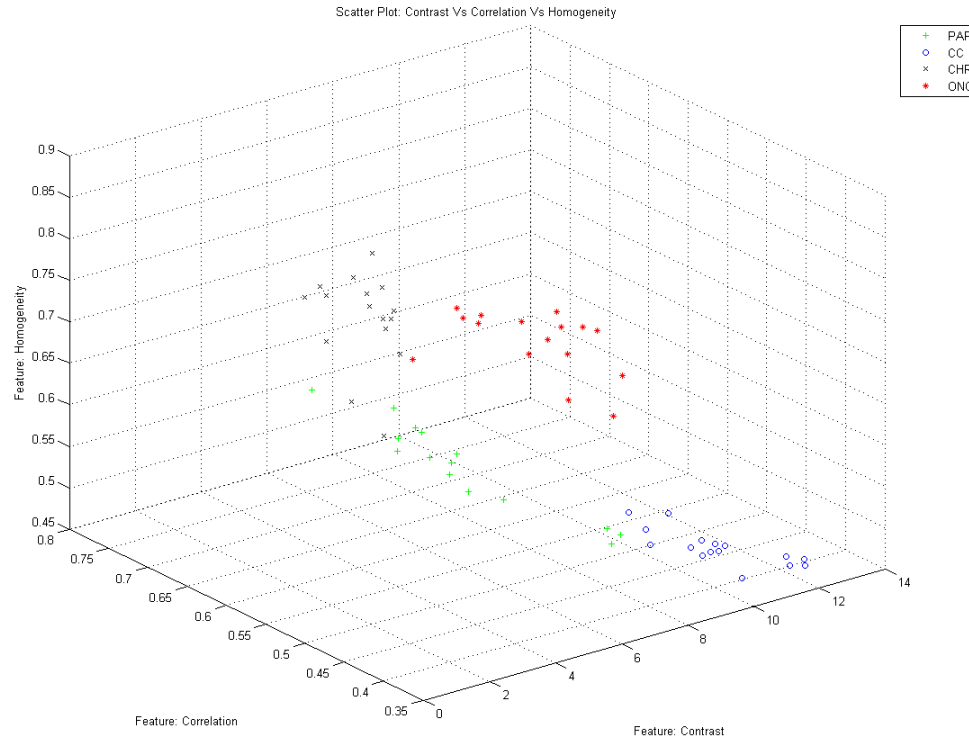
usually found inside blood vessels which tend to appear white under the microscope. This information leads to the identification and separation of RBCs from other objects with in the image. Again the knowledge of the system guides the selection of the regions over which these properties were calculated which is important in this case. These features are calculated using a co-occurrence matrix [86]. The textural features are calculated for the entire image, sub-image and the individual cellular structures. Table 1.1 lists the values of some of the features extracted. The same tool will be able to quantify features from other form of cancer imaging data by extracting the features selected by the user and over a particular ROI.

Once the features of interest are extracted they were passed on to the classification system which decided the subtype of cancer image based on the feature values.

#### **4.1.5 Classification**

With various properties of the tissues extracted and quantified, their variations from a particular value well help us classify the different subtypes correctly. If one studies Table 1.1 closely he would find that a single feature alone couldn't differentiate between the different subtypes. This derives directly from the fact that that renal tumor subtypes exhibit several common morphological characteristics. Adding additional features to the analysis however, (see Figure 4.2) help separate the classes from each other spatially. It is evident from the scatter plot that these should be easily separated and classified if a correct combination of features extracted and used as is the case. To attempt the most robust classifier 8 knowledge-based features selected for extraction by the pathologist are used by the classifier. The classification scheme used here implements a simple, multi-class Bayesian decision rule that assumes multivariate Gaussian distributions for the data. To estimate the ability of the classification rule to correctly predict the class of a new, unknown sample, complete leave-one-out cross-validation was

performed. Results showed that the classifier correctly predicted the class of an unidentified sample 98.4%, or 63/64 of the time.



**Figure 4.2** Scatter Plot showing distribution of Images for three co-occurrence features: Contrast, Correlation and Homogeneity (+, Papillary; o, Clear Cell; x, Chromophobe; \*, Oncocytoma)

The classification results prove that the strategy of feature extraction based on expert pathologist knowledge results in high classification accuracy (98.4%). This classification accuracy is particularly high due to the fact that the user oversees the whole process from image processing, feature extraction and quantification to classification. Although these features correspond to visual properties traditionally used by the pathologist and usually correspond to properties not common to all subtypes but same set of features will not be useful for different subtypes of cancers and hence a pathologist expert input is necessary. In the present case the presence of lipid structures, Red blood

cells and blood vessels varies among different subtypes of renal cancer and their presence was quantified using their textural properties.

## **4.2 Morphological, textural and wavelets based features**

In this section, we present an image quantification and classification method for improved pathological diagnosis of human renal cell carcinoma (RCC). This method combines different feature extraction methodologies, and is designed to provide consistent clinical results even in the presence of tissue structural heterogeneities and data acquisition variations. The methodologies used for feature extraction include image morphological analysis, wavelet analysis and texture analysis, which are combined to develop a robust classification system based on a simple Bayesian classifier. We have achieved classification accuracies of about 90% with this heterogeneous dataset. The misclassified images are significantly different from the rest of images in their class and therefore cannot be attributed to weakness in the classification system.

### **4.2.1 Feature Extraction**

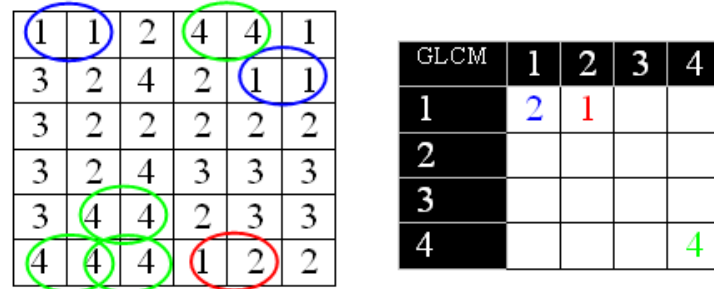
The input images after initial processing (section 4.1) are segmented in four-level grayscale image. We then analyze images from different RCC subclasses and try to predict the variations in these regions (e.g., the gray level one represents the nuclei in the images). The size of nuclei in one subclass may be different from another subclass, and can be used as one of the differentiating features between the classes. In practice, we may find a feature that differentiates between two subclasses, but this feature may not be useful for other subclasses. This necessitates finding a larger set of features that can differentiate between more subclasses. For this purpose, we have used different methodologies and have combined their results into a set of significant features, which are then used to improve the accuracy of our classification system.

The first method uses gray-level co-occurrence matrix (GLCM), also known as the gray-level spatial dependence matrix [87, 88]. Unlike the texture filters, which

provide a statistical view of texture based on the image histogram without providing any information about shape, the GLCM method combines textures and morphological statistics into one matrix.

The GLCM is computed by calculating how often a pixel with the intensity (gray-level) value  $i$  occurs in a specific spatial relationship to a pixel with the value  $j$ . Each element  $(i, j)$  in the resulting GLCM is the sum of the number of times that the pixel with value  $i$  occurs in the specified spatial relationship to a pixel with value  $j$  in the input image.

GLCM computation on our four-level grayscale images generates a four by four matrix, and an example is shown in Figure 4.3. Figure 4.3(a) represents a portion of a 4-level gray scale image with elements (1,1), (1,2) and (4,4) indicated for co-occurrence of immediate horizontal neighbors, using an offset mask of  $\begin{bmatrix} 1 & 1 \end{bmatrix}$ . Figure 4.3(b) shows the corresponding entries in the GLCM using the sum of highlighted elements.



**Figure 4.3** GLCM computation using 4-level grayscale images. (a) Representation of 4 level grayscale image (b) GLCM for highlighted elements in image (a)

The gray-level co-occurrence matrix can reveal certain properties about the spatial distribution of the gray levels in the image. For example, if the entries in the GLCM diagonal data are large, the regions are contiguous and the texture is coarse. With a small offset and the large concentrated entries, each diagonal element represents an image area of the corresponding gray-level region of interest. In our implementation, gray-level '1' represents the nuclei, so the GLCM element (1,1) shows the count of total nuclei area in the image in terms of pixels. This count divided by the image size (1200x1600) gives the



normalized nuclei area in the image and is used as one of our desired features. A few significant features extracted from GLCM are shown in Table 4.2.

Table 4.2 Features extracted from GLCM for clear cell (CC), chromophobe (CHR), renal oncocytoma (ONC) and papillary (PAP) with mean and standard deviations.

	CC	CHR	ONC	PAP
<b>GLCM (1,1)</b>	0.1095±0.0323	0.0608± 0.0105	0.0918± 0.0289	0.1577± 0.0347
<b>GLCM (1,2)</b>	0.0321± 0.0079	0.0194± 0.0037	0.0285± 0.0083	0.0451± 0.0083
<b>GLCM (2,2)</b>	0.3988± 0.0478	0.5505± 0.0628	0.4687± 0.1124	0.3614± 0.0587
<b>GLCM (2,4)</b>	0.0577± 0.0141	0.0547± 0.0057	0.0556± 0.0242	0.0427± 0.0106
<b>GLCM (3,3)</b>	0.0311± 0.0459	0.0057± 0.0030	0.0217± 0.0465	0.0029± 0.0039
<b>GLCM (4,4)</b>	0.2613± 0.0922	0.2260± 0.0763	0.2309± 0.0848	0.2943± 0.0637

The same GLCMs can be used to derive other statistics about the texture of an image. The most commonly used statistics include:

**Contrast** – it measures local variations in the gray-level co-occurrence matrix.

$$\sum_{i,j=0}^{N-1} P_{i,j} (i - j)^2$$

**Correlation** – it measures the joint probability occurrence of the specified pixel pairs.

$$\sum_{i,j=0}^{N-1} P_{i,j} \frac{(i - \mu_i)(i - \mu_j)}{\sqrt{\sigma_i^2 \sigma_j^2}}$$

**Energy** – it is also known as uniformity or the angular second moment, and provides the sum of squared elements in the GLCM.

$$\sum_{i,j=0}^{N-1} P_{i,j}^2$$

**Homogeneity** - it measures the closeness of the distribution of elements in the GLCM to the GLCM diagonal elements.

$$\sum_{i,j=0}^{N-1} \frac{P_{i,j}}{1 + (i - j)^2}$$

**Entropy** - it measures the randomness between the elements of GLCM

$$\sum_{i,j=0}^{N-1} P_{i,j} (-\ln P_{i,j})$$

These textural statistics are computed, and their results are listed in Table 4.3.

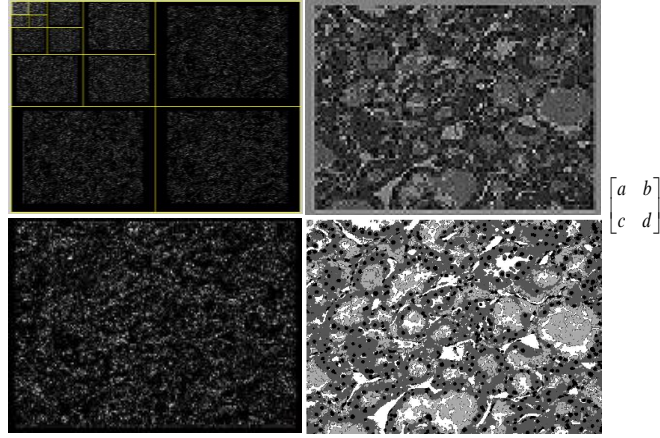
Table 4.3 Statistical Features extracted for clear cell (CC), chromophobe (CHR), renal oncocytoma (ONC) and papillary (PAP) with mean and standard deviations.

	CC	CHR	ONC	PAP
<b>Contrast</b>	0.5712± 0.0726	0.4937± 0.0492	0.5398± 0.2004	0.4592± 0.0893
<b>Correlation</b>	0.7496± 0.0173	0.7340± 0.0610	0.7392± 0.0701	0.8252± 0.0270
<b>Energy</b>	0.2640± 0.0283	0.3759± 0.0366	0.3140± 0.0778	0.2596± 0.0193
<b>Homogeneity</b>	0.8811± 0.0069	0.9040± 0.0091	0.8885± 0.0313	0.8944± 0.0148
<b>Entropy</b>	2.3825± 0.1960	1.9577± 0.0787	2.2077± 0.2663	2.3110± 0.0639

The use of wavelet transform [86] can also improve feature extraction by performing multi-resolutions analysis of the image. Wavelets are mathematical functions that decompose data into different frequency components, and then study each component with a resolution matched to its scale. Because of its representation of piecewise-smooth signals and fractal behavior owing to its multi-resolution, this method has been successfully used for many biomedical imaging applications [89, 90].

We have used our four-level grayscale segmented images for wavelet analysis. Bi-orthogonal wavelet pairs of the third order, a family of B-Splines, are used as the wavelet basis. The transformation generates four component sub images, known as Approximation and Detail (Horizontal, Vertical and Diagonal). Figure 4.4 shows the 4-level DWT image, selected sub-images, and the original grayscale image.

The Wavelet coefficients from the previous stage are processed to enhance the objects of interests. These coefficients contain positive and negative intensities. In post processing, we take the absolute of these intensities and reduce the number of gray levels to four. The remaining analysis of the wavelet images is identical to the process above for four level grayscale images. Each sub image is analyzed using GLCM as well as textural analysis of GLCM by finding properties like contrast, correlation, etc.



**Figure 4.4** (a) 4 level DWT image, (b) Level three approximation component, (c) Level two horizontal detail component, (d) 4-level grayscale ONC image.

Features are computed from 4 levels of wavelet transform. Every sub image from every level contributes to a large cumulative set of features. These features are ranked based on how well they can discriminate between the RCC images in the training database. Some significant features obtained through the wavelet analysis are listed in Table 4.4. Once the features of interest are extracted, they are used by the classification system to determine the subtype of cancer image.

**Table 4.4** Features extracted after DWT for clear cell (CC), chromophobe (CHR), renal oncocytoma (ONC) and papillary (PAP) with mean and standard deviations

	CC	CHR	ONC	PAP
<b>GLCM(1,1) (level1-Approx)</b>	0.1072±0.0297	0.0582±0.0109	0.0918±0.0265	0.1601±0.0404
<b>GLCM(1,1) (level1-Diagonal)</b>	0.6327±0.0261	0.6951±0.0383	0.6844±0.0882	0.6915±0.0370
<b>GLCM(1,2) (level1-Diagonal)</b>	0.1122±0.0051	0.0920±0.0089	0.0986±0.0237	0.0988±0.0110
<b>GLCM(2,2)(level1-Horizontal)</b>	0.1260±0.0103	0.1053±0.0116	0.1137±0.0296	0.1080±0.0146
<b>Homogeneity (level 1- Vertical)</b>	0.7618±0.0212	0.7976±0.0222	0.7881±0.0629	0.8048±0.0317
<b>Energy (level 1-Diagonal)</b>	0.4054±0.0329	0.4874±0.0528	0.4627±0.1131	0.4828±0.0513
<b>Entropy (level 1-Horizontal)</b>	0.3941±0.0364	0.4628±0.0476	0.4429±0.1130	0.4677±0.0555
<b>Energy (level 2-Diagonal)</b>	0.3221±0.0420	0.3896±0.0534	0.3652±0.1157	0.3969±0.0562

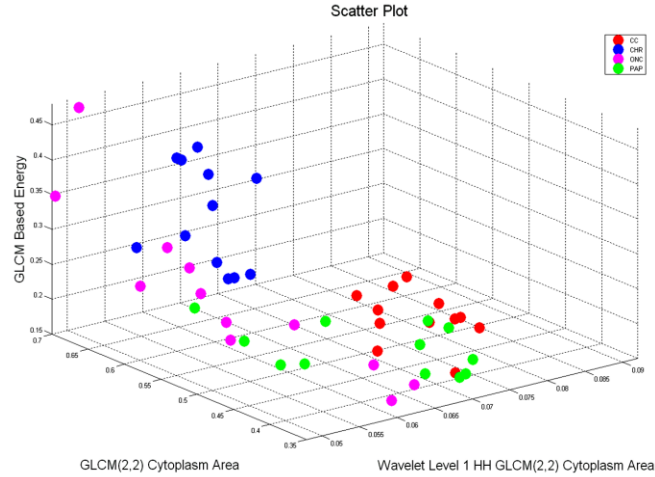
### 4.2.2 Classification

Various features extracted from the RCC tissue images using different methodologies are used to classify the subtypes correctly. The analysis of data in Tables 1, 2, and 3 shows that the RCC subtypes have considerable similarities that make it difficult for a single feature to differentiate all the subclasses. This problem can be addressed by increasing the dimensionality and adding more features. As shown by the scatter plot in Figure 4.5, the different subclasses can be separated by using three features.

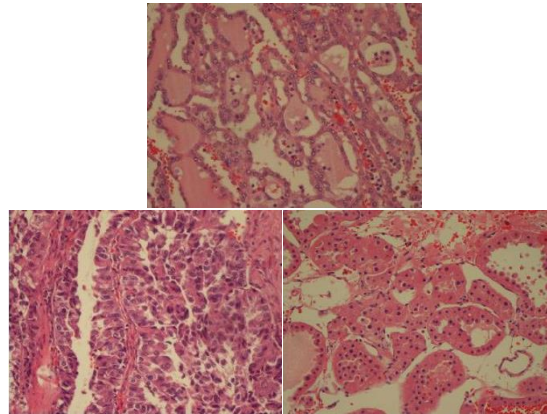
We have used simple, multi-class Bayes classifier assuming multivariate Gaussian distributions to predict RCC image subclasses. The leave-one-out cross-validation method is used to evaluate the ability of our features and the classifier to correctly predict unknown images. Our results show that by using the features listed in Table 4.5, our classifier can correctly predict the class of an unidentified sample with an accuracy of 87.5 % from our significantly heterogeneous image data. It is interesting to note that most of the misclassified images are significantly different from other images in their own class, thereby contributing to false detection.

Table 4.5 List of features selected for the best classification performance with mean and standard deviations.

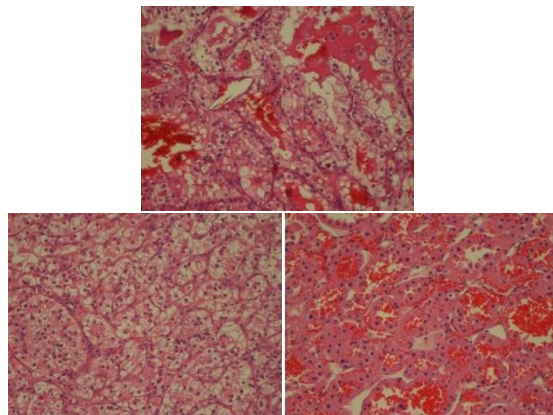
	<b>CC</b>	<b>CHR</b>	<b>ONC</b>	<b>PAP</b>
<b>GLCM(1,1) (level1-Approx)</b>	0.107±0.0297	0.058±0.0109	0.091±0.0265	0.160±0.0404
<b>GLCM(2,2)(level1 Diagonal)</b>	0.067±0.0031	0.054±0.0060	0.058±0.0200	0.057±0.0083
<b>GLCM(1,1)</b>	0.109±0.0323	0.060±0.0105	0.091±0.0289	0.157±0.0347
<b>GLCM(2,2)</b>	0.398±0.0478	0.550±0.0628	0.468±0.1124	0.361±0.0587
<b>Homogeneity</b>	0.881±0.0069	0.904±0.0091	0.888±0.0313	0.894±0.0148
<b>Energy</b>	0.264±0.0283	0.375±0.0366	0.314±0.0778	0.259±0.0193



**Figure 4.5:** Scatter plot showing distribution of images for three features: GLCM based Energy, GLCM based Diagonal component representing cytoplasm area and Wavelet level 1 Diagonal detail GLCM component (2,2)



**Figure 4.6:** (Center) PAP image misclassified as ONC. (Left) another PAP image (Right) ONC image



**Figure 4.7:** (Center) CC image misclassified as ONC. (Left) another CC image (Right) ONC image

Due to significant difference of these images from other images in their class, some images are misclassified but even the misclassified images are unanimously annotated to one class by most of the features. Figure 4.6 and 4.7 shows two examples of the misclassified images along with an image from their correct class as well as an image to the incorrect class to which the image was annotated.

Although this study showed promising results, much work is still needed towards the eventual goal of a clinical image classification system for routine pathological use.

Image features from a single methodology may be good enough for a less varying dataset. But as we increase the complexity and heterogeneity of the images, it becomes difficult to find a set of features that can consistently produce accurate results. By combining features obtained using different methodologies, we show that the feature set becomes more robust and can achieve accurate and consistent results. In addition to improving the feature extraction and classification process, the standardization of tissue sample preparation and the image acquisition process are also important factors. Further, in actual practice the pathologist only concentrates on part of the image and bases his or her classification on the specific region of interest. Manual segmentation of ROI considerably improves the classification accuracy approaching 100% for some datasets. Automatic segmentation of these areas is by itself a problem of considerable complexity. It is thus important to combine these two problems into one system and use areas of high correlation with the training database. Work is ongoing to provide a sophisticated system to assist pathologists in their diagnosis and early cancer detection leading to improved survivability of renal cancer patients.

### **4.3 Cellular features of elliptical models of segmented nuclei clusters**

In this section, we present the results of our computer aided diagnosis system for subtype classification of Renal Cell Carcinoma pathological images based on features of

individual cells. Traditionally, cancer diagnosis is done by an expert pathologist by studying biopsy tissue under a microscope. Heterogeneity in patient tissue samples, variation in sample preparation process, intra and inter-observer variability make the classification task quite complex and challenging. This requires use of a computational diagnosis system to improve the repeatability and accuracy of the process as well as assist the pathologist in decision making. Previous sections proposed different methodologies and used a combination of several features, derived from the complete image, as a solution to this problem. However, pathologists inherently consider only a part of the biopsy slide and also take into consideration the features of the individual cells. To replicate this human behavior, we need to segment individual cells in the tissue images to extract features for classification. This, in-turn, poses a significant challenge as some images have dense nuclei clusters which are very difficult to segment. Based on our recent work, we used concavity based ellipse fitting technique to segment the nuclei clusters and then determine individual cell features. We report high classification accuracy (94%) on a heterogeneous tissue image data set which has significant intra-class variations. We also hope that this methodology will help pathologist's decision making in the clinical setting.

We obtained reasonably good classification accuracy in our previous works[49, 54, 79]. In[49, 54], we extracted textural and wavelet based features from the whole image. In [79], we excluded necrotic regions and large lumen spaces to guide our technique closer to a pathologist's method of analyzing only the relevant regions in a biopsy tissue image. In this work, we further extend our technique to imitate yet another crucial step followed by pathologists. Pathologists not only look at the significant portions of an image, but they also look at the characteristics of individual nuclei. Our present work involves feature extraction from individual nuclei and thus encapsulates several inherent procedures followed in manual grading by pathologists. Since necrotic regions and large lumen regions lack nuclei, they are automatically rejected for feature

extraction. We also extract features from regions around the nuclei. The characteristics of these regions are an important criterion to distinguish the subtypes. For example, CH subtype is clearly distinguished by a clear halo around the nuclei while in PA, nuclei are clustered together and the region around a single nucleus may be overlapped by another nucleus. ON and CC nuclei are far apart, however, the region around their nuclei are eosinophilic and clear respectively. Our method combines relevant region analysis with knowledge-based feature extraction. Our automated classification system classifies the renal cell carcinoma images into four subtypes with minimal user interaction and reasonable accuracy showing the potential for future clinical use.

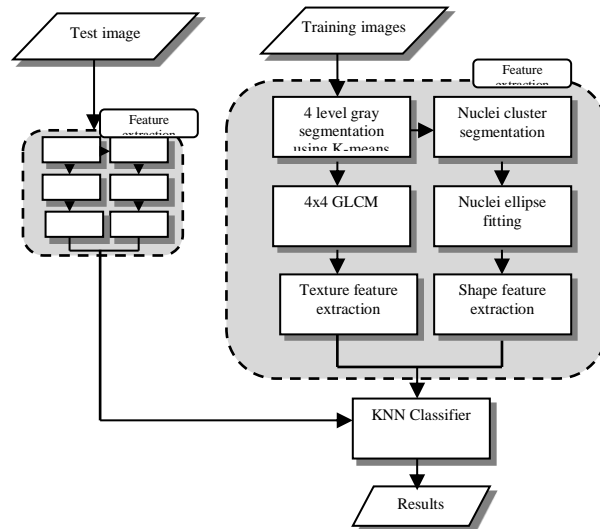
In [49], we used knowledge-based RCC features to obtain high classification accuracy using our test data set which was selected by the pathologist as a good representative of each RCC subtype. However, in practice, clinical image data is highly heterogeneous with significant variations in the images of each RCC subtype. Our algorithm in [49] gives reduced classification accuracy when used for significantly heterogeneous images within each subtype class. In addition to the heterogeneity, the tissue samples also contained necrotic regions which contribute to the reduced accuracy. In [54], we designed a new methodology to overcome the reduced accuracy in the presence of heterogeneous data. We extracted features using a combination of morphological analysis, wavelet analysis and texture analysis. In [79], we augmented our knowledge-based classification system with automatic region of interest (ROI) selection and rejecting necrotic zones. We demonstrated the importance of intelligent ROI selection to reduce computation time and increase classification accuracy. Taking a step further and imitating pathologist practice of basing decision on individual cells, we propose a new system which segments the nuclei clusters based on ‘concavity based ellipse fitting methodology [91]. We compute individual cell features and then use a KNN classifier to finally classify RCC tissue image to one of its four sub-classes. We focus on nuclear segmentation and derive features for each segmented nuclei with several



angle and distance measures. We report high classification accuracy (94%) on a heterogeneous tissue image data set which has significant intra-class variations. We also hope that our proposed methodology will help pathologists in their decision making process.

#### 4.3.1 Methodology

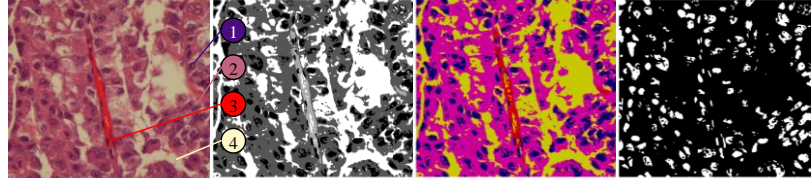
We use standard photo micrographs of hematoxylin & eosin (H & E) stained biopsy tissue sections as our image dataset. We process these images using the methodology shown in Figure 4.8. First, we color segment each image and then convert it into four-level grayscale images (one level for each color corresponding to nuclei, cytoplasm, red blood cells and unstained clear tissue). Next, we segment the nuclei using a mask corresponding to nuclei derived from color segmentation. Nuclei shape features are extracted using the fitted ellipse model over segmented nuclei while the texture features inside and outside nuclei are computed using GLCM. These features were then used to train the KNN classifier with subsequent classification of the unknown images into subtypes of the RCC. We will describe the detailed processing steps in the remainder of this section.



**Figure 4.8** Flowchart for the overall methodology

### 4.3.2 Image acquisition & image color segmentation

The image acquisition and initial processing of the images was done in a similar fashion as discussed in section (4.1) above. The H&E staining with red blood cells and the background results in four distinct colors in the images. The color and intensity of the images, however varies considerably due to variations in sample preparation and the image acquisition processes. Consequently, to be consistent with tissue staining, we segment the images into four-level grayscale images, each level corresponding to a mask for one out of four object categories, that is, nuclei, gland, cytoplasm and red blood cells. A large variation in intra-sample color and intensity requires some intelligent processing to segment the RGB images into quantized grayscale images representing region masks. We use K-means clustering for our RCC subtype images. We start with the fixed initial values of the staining colors as the means of the  $k=4$  clusters. The K-means algorithm adjusts to the variation in the images by shifting the cluster means and updating the pixel assignments. Figure 4.9 shows the results of this segmentation. Our seeded K-means is constrained and can only shift means within a specific range thereby capturing stain variations but cannot shift means enough to change color classes beyond standard staining colors. We also provide visual feedback to the user to validate the color segmentation and provide slider control if the user wants to bias the color classes. All these methods leave a little room for major segmentation discrepancies while our method is robust enough for minor discrepancies in segmentation.



**Figure 4.9** K-Means segmentation (Left to right): 1) original image; 2) gray level segmented image; 3) segmented pseudo color image; 4) Nuclei mask

#### 4.3.3 Nuceli cluster segmentation and ellipse fitting

Nuclei mask generally contains clusters of nuclei which need to be split, before being considered for computation of their individual features. Classical splitting-nuclei methods such as watershed can only split simple clusters with small overlap. These algorithms frequently produce erroneous results which can seriously affect overall image statistics and can lead to misclassification. We use a methodology for cluster



**Figure 4.10** (Left) Cluster segmentation by watershed method (Right) Cluster segmentation by notch detection and ellipse fitting method

segmentation using notch detection and ellipse fitting proposed in our previous work [91] to achieve this goal. A comparison of watershed segmentation with our algorithm is shown in Figure 4.10.

The segmentation process involves certain steps as explained below:

a) *Preprocessing*: The nuclei masks are obtained during image color segmentation process as explained in section 2.2 above. These masks suffer from salt and pepper noise and small holes in the nuclei area due to misclassification during K-means. These artifacts are removed using standard morphological

techniques like image closing and opening. The cleaned up masks are then processed for concavity detection.

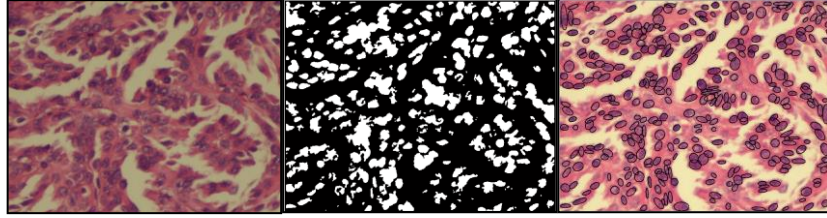
*b) Notch detection:* Each cluster in the nuclei masks is processed for detection of sharp notches along the boundary of the cluster. These notches are considered as possible points for cluster splitting. The notches on the surface can be either concave or convex. The points where two cells touch each other are concave. We detect this by dividing the complete cluster boundary into small linear segments and determining the tangential vectors for every segment using endpoints. If  $p_1, p_2, p_3$  are end points of two consecutive segments then their tangential vectors  $\mathbf{a}$  and  $\mathbf{b}$  are given by equation 7. Considering each pair of consecutive segments as vectors, we compute their cross product. As  $\mathbf{a}$  and  $\mathbf{b}$  has zero z-component, cross product will be in z-direction. The sign of cross product determines the surface to be convex or concave while the angle between the vectors determines the sharpness of these notches. Equation 7 explains the process.

$$\mathbf{a} = [p_{2x} - p_{1x}, p_{2y} - p_{1y}, 0] ; \mathbf{b} = [p_{3x} - p_{2x}, p_{3y} - p_{2y}, 0] ; \sin \theta = \frac{1}{|\mathbf{a}| |\mathbf{b}|} [\mathbf{a} \cdot \mathbf{b}_y - \mathbf{a}_y \mathbf{b}_x] \quad (7)$$

For concave portion of an edge,  $\sin \theta > 0$  and for convex portions,  $\sin \theta < 0$ . Notches can be detected by using low positive threshold. Multiple concavities may be detected in near vicinity of each other. The concavity with maximum angle change amongst its neighbors is selected as a seed point for cluster splitting.

*c) Straight line segmentation:* After identifying the possible seed points for cluster splitting, an iterative process of cluster splitting is used where a split is created based on two nearest seed points and then each sub-cluster is evaluated for further split. This initial segmentation of the clusters is further refined during the ellipse fitting session.

*d) Ellipse fitting:* We use direct ellipse fitting method proposed by Fitzgibbon et al. [81] for its accuracy and performance. Our ellipse fitting process iteratively fits an ellipse to the boundary points of each portion of the segmented cluster. At every step, we check the overlap of the present ellipse with the previously fitted ellipses to avoid over segmentation. Figure 4.11 shows the nuclei segmentation with ellipse fitting process.



**Figure 4.11** Nuclei cluster segmentation for PA RCC tissue image. (a) RGB image, (b) Binary nuclei mask, (c) Individual nuclei marked on RGB tissue image using ellipse fitting.

#### 4.3.4 Feature extraction and selection

We use shape and texture features of the individual cells. The shape features are purely dependent on the elliptical model of each nucleus while the texture features are based on the GLCM within nuclei and immediate neighborhood of the nuclei (e-glcm -> GLCM of exterior region). Our previous work [54] show independent results primarily from GLCM alone provide good accuracy. The reason for using elliptical descriptors is that they are closer to the human interpretation e.g. how round are the nuclei? How big are the nuclei clusters? What is average nuclei count in each cluster? In addition these descriptors are more useful for the next CAD step of cancer grading which is primarily done on the basis of morphological shape analysis. These features are explained below:

##### Shape based features:

After ellipse fitting, for each nucleus we have basic ellipse parameters like major axis, minor axis, eccentricity and angle between the major axis and x-axis. Using these

parameters we can determine shape features of the individual nuclei and also their distribution in the image. The features we used for our classification system are listed below.

1. Nuclear area given by

$$A = \pi \times a \times b \quad (8)$$

where a is major axis and b is minor axis.

2. Nuclear Eccentricity given by

$$E = \sqrt{1 - \frac{b^2}{a^2}} \quad (9)$$

where a is major axis and b is minor axis.

3. Inter-nuclear distance in neighborhood nuclei,  $N_{ij}$ : We select a fixed number of neighbors around each nucleus and calculate its distance from each neighbor. An average value of the inter-nuclear distance is used as a feature.

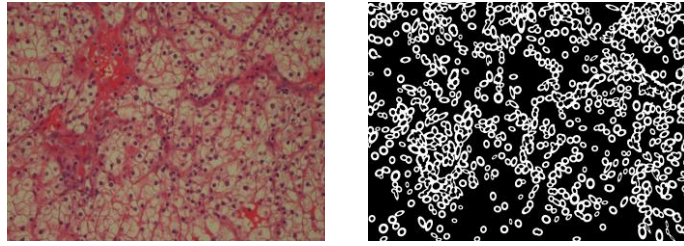
4. Major axis alignment with neighboring nuclei: In PA subtype, sometimes there are chains of cells oriented in the same direction. Also, in compact tumor regions, cells are pressed against one another. Rapidly dividing cells might also result in a particular orientation. To capture these effects, we calculate the relative alignment of major axis of a given nucleus with its neighbors.

5. Nuclei density: Since the image size is constant for all images, nuclear density provides an implicit measure of the degree of compactness of tumor region.

#### Texture features:

Our texture feature extraction is based on gray level co-occurrence matrix (GLCM) for each individual cell. The area inside each nucleus as well as area outside each nucleus is considered separately. The outside region is a similar elliptical region with both major and minor radii double than the radii of nucleus under consideration as shown in Figure 4.12. The GLCM computed within each area captures the frequency that

a gray-level value occurs adjacent to another gray-level value [92]. As we already segmented the images into four gray level intensities, our GLCM is a  $4 \times 4$  matrix  $glcm(i,j)$ . One GLCM matrix of size  $4 \times 4$  represents one spatial relation (e.g. horizontal) between the intensities of the image. Therefore, we calculated four GLCM matrices to cover all the four spatial relations (horizontal, vertical, diagonal at the angle of 45 degree and diagonal at -45 degrees) between the intensities from the image and take the average of these four to present the overall spatial relation of the gray level intensities within the image.



**Figure 4.12** (Left) CC input image. (Right) Mask of regions outside the nuclei.

Some of the GLCM measurements can be easily correlated to the tissue properties. For example, with color selection as shown in Figure 4.9,  $glcm(1,1)$  is a measure of amount of nuclear stain in the region under consideration,  $glcm(2,2)$  is a measure of amount of cytoplasmic stain,  $glcm(4,4)$  is a measure of amount of unstained tissue and  $glcm(1,2)$  combined with  $glcm(2,1)$  is a measure of amount of edges between nuclei and cytoplasm regions.

In our four color selection, color 3 corresponds to the red blood cells (RBC). Presence of RBC cannot be attributed as one of the properties to any RCC sub-class and therefore it is expected that all the GLCM with color 3 that is  $glcm(3,1)$ ,  $(3,2)$ ,  $(3,3)$ ,  $(3,4)$ ,  $(1,3)$ ,  $(2,3)$  and  $(4,3)$  will not contribute as a significant feature. In addition, the GLCM matrices obtained are near symmetric, so we included average of upper and lower triangle of GLCM matrices as part of our features set. Each GLCM matrix triangle has 10 elements and eliminating 3 RBC related elements reduces the set to 7 GLCM based

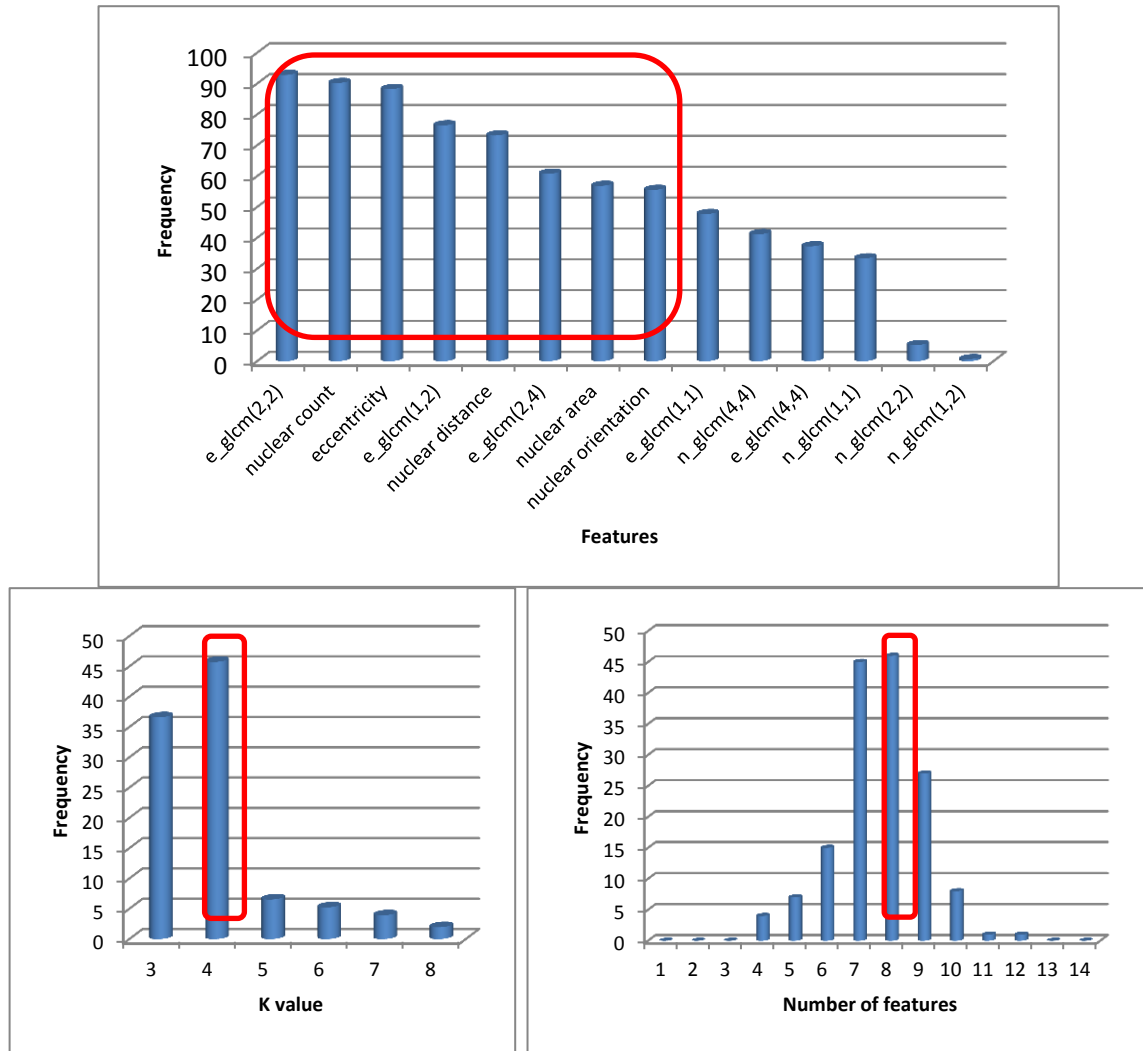
features. Our total feature space includes 7 GLCM features inside the nuclei ellipse, 7 GLCM features in the area around nuclei ellipse and 5 shape based features.

#### **4.3.5 Classification**

We used the features extracted in the previous step for classification of images into the subtypes using the KNN classifier. The KNN classifier was selected based on the fact that our image data set have significant intra-class variation. Each test image is expected to have some similar images in its class, although it may not be similar to all the images in the corresponding class.

Our data consists of 4 subtype classes with 12 images in each class. We split the data in three folds with 32 images used for training and 12 images were used as test data. We performed leave-one-out cross validation (CV) on training images varying all combinations of features and values of K. The selected models were used to perform external validation (EV) using the test images. Multiple models have same CV accuracy. Figure 10 shows a relationship of these CV models with their corresponding EV models. It can be seen that CV and EV results are well correlated with accuracies for best models above 90%. Figure 8b shows how frequently a value of K appeared for best model. Analyzing the top models, with accuracy greater than 90%, we pick our best feature size as 8(Figure 8c). We rank our features (Figure 8a) based on how frequently they appear in top models. Figure 9 shows a bar graph with averages values of these features highlighting significant variation in these features for different subtypes.





**Figure 4.13** Statistics for top models. (a -top) Most significant features based on how frequently they showed up in top models (b - bottom left) Best K-value selection for KNN. (c -bottom right) Number of features used by top models

Table 4.6 Confusion Matrix

	CC	CH	ON	PA
CC	91.7%	8.3%	0%	0%
CH	0%	100%	0%	0%
ON	0%	8.3%	75%	16.7%
PA	0%	8.3%	0%	91.7%

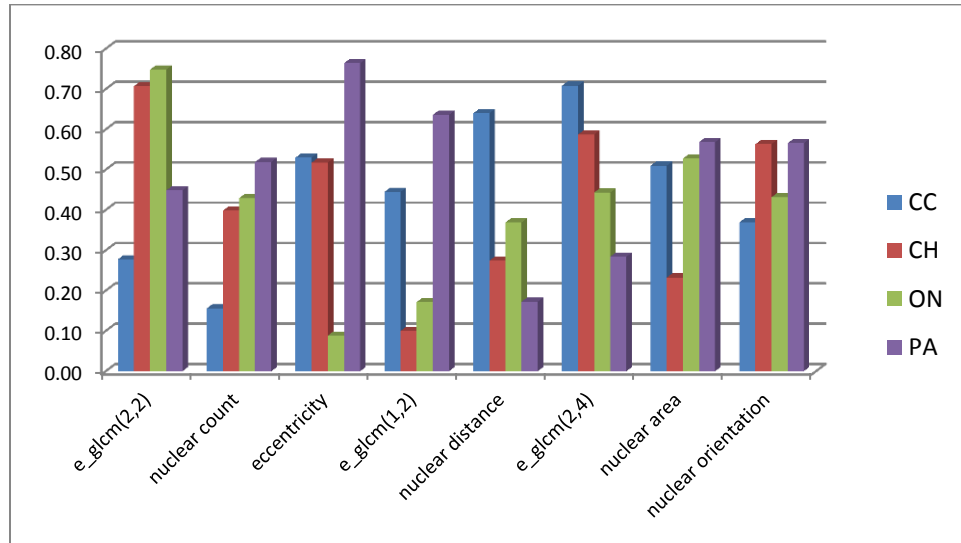


Figure 4.14 Distribution of normalized features for different subtypes

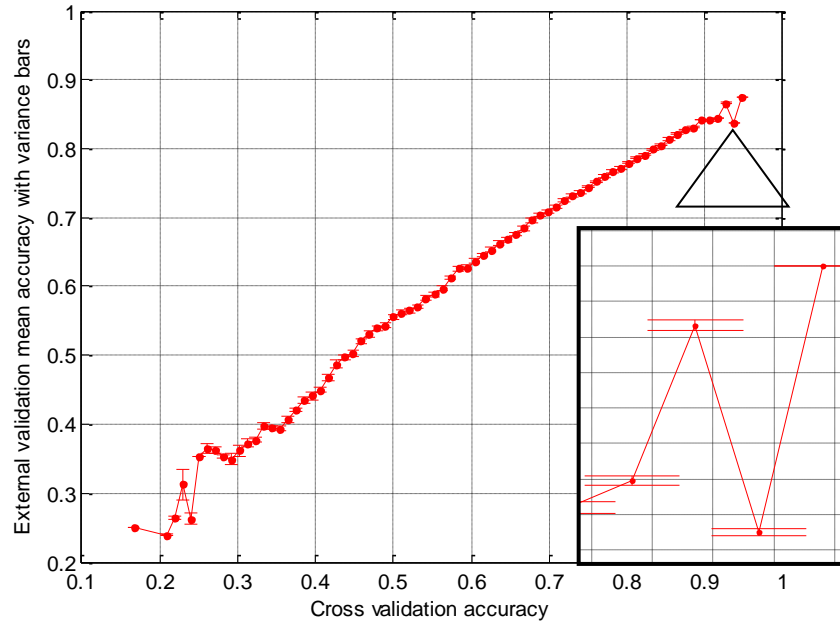


Figure 4.15 Correlation between cross validation and external validation results

### 4.3.6 Results and discussion

Our results show that each subtype of H&E stained RCC images can be classified with high accuracy. We obtained above 90% classification accuracy for our best models. The confusion matrix in Table 4.6 also supports that the error is low and is distributed. Even the subtypes like CH and ON which are generally hard to classify are identified correctly. Although, the rate of occurrence of some sub-types is much larger than the others, correct diagnosis of rarely occurring subtypes is also critical. If a malignant class is misclassified as the benign class ON, the patient may be left un-treated resulting in serious future consequences. This fact further supports the use of our system being used as a decision support in aid of pathologist to verify and validate their diagnosis.

Selected features of our model (Figure 4.13) highly correlate with the human knowledge of feature selection and the explanation of some of these is given below.

***Inter-nuclear distance:*** PA images have a property of forming dense nuclei clusters and it is expected that  $N_{ij}$  computed from the near neighbors will be less in case of PA images as compared to the other classes like ON where the nuclei are well separated. Our statistical results (Figure 4.14) show that PA images have least inter-nuclear distance.

***Nuclear eccentricity:*** ON class has round nuclei which correspond to lower eccentricity values. This fact is also well supported by the computed statistics (Figure 4.14).

***e-glc<sub>m</sub>(2,2):*** This measure has been found lowest in CC class due to presence of clear tissue around nuclei, relatively lower in case of PA due to touching nuclei and relatively larger for CH and ON classes (Figure 4.14)

***Nuclear count:*** PA images have large nuclei clusters and overall larger nuclear density in comparison to the other subtypes. The fact is accordingly observed with PA class having highest and CC with lowest nuclear count (Figure 4.14).

Our proposed methodology takes care of the unwanted necrotic regions automatically by considering only area inside nuclei and in immediate vicinity around nuclei. Figure 14 shows some larger black zones which are devoid of nuclei and were not used for computation of statistics thereby maintaining the higher accuracy proposed in our previous work[54].

With our improved classification system replicating a pathologist's method, we are highly motivated to apply this system for classification of other cancer types. Our results are consistent with our previous work [79] which showed improved classification accuracy by using better ROI segmentation. By using nuclei segmentation, only selective zones within the image contribute to the statistical computation, thereby resulting in much improved and robust system. Evaluation of nuclear features is also expected to provide important features for computer assisted grading of RCC subtypes. For example, CC grade III and IV are distinguished by nuclei and nucleoli size variation along with texture of nucleoli[5]. Our methodology, thus can also be acclimated to perform computer assisted grading of various cancer types.

## **CHAPTER - V**

### **RCC GRADING BASED ON FUHRMAN NUCLEAR GRADE**

In 1932, Hand and Broders study [93] found that Renal Cell Carcinoma (RCC) grade was associated with the outcome of the patient. Patients with higher grade carcinoma had more mortality and shorter survival time than patients with lower grade. Since then different grading systems have emerged based on a number of studies which examined these relationships. Most RCC grading systems [94] target combination of nuclear and nucleolar characteristics, and a few have also analyzed the cell type and tumor architecture. The grading systems [95-97] have shown prognostic merit of the RCC grading and its correlation with the likelihood of metastases or local recurrence. Although RCC grading correlates with survival and stage, it is not an independent, significant predictor. Studies like [98] have shown other factors like proliferating cell nuclear antigen (PCNA) expression as a prognostic indicator for RCC.

There is a consensus with regard to the utility of RCC grading, but there is no agreement regarding which grading system should be used. Each system has significant advantages and disadvantages.

#### **5.1 Fuhrman Grading**

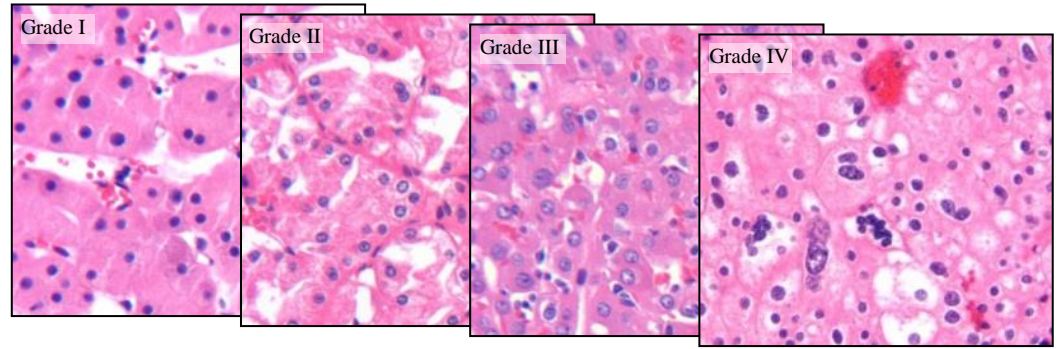
Most pathologists in North America use the Fuhrman et al. grading system[99], and many studies that have examined the utility of RCC grading have employed this system. It is a four-tiered system where grade I carries the best prognosis and grade IV the worst. The major criterion that distinguishes each tumor grade is the presence of a nucleolus, its size, and the magnification at which it can be observed. Based on the fact that the system is based on just the appearance of the nuclei of the cancer cells, rather

than the appearance or structure of the cells as a whole, it is also termed as Fuhrman Nuclear Grade.

Nuclear characteristics used in the Fuhrman Grade particularly indicate how actively the cells are making protein. These characteristics include size and shape of the nucleus as a whole, number and size of nucleoli (Nucleoli are organelles found in the cell nucleus which make ribosomes which in turn are protein making factories. More nucleoli implies more active protein synthesis) and chromatin clumping. Chromatin is the substance of chromosomes, which includes DNA, chromosomal proteins, and chromosomal RNA. Chromatin stains strongly with basic dyes. It is thought that the chromatin is most deeply stained when it is most condensed and inactive. Well differentiated tumors are recognized as exhibiting orderly stratification, obvious cellular bridges, and keratin pearl formation. In contrast, poorly differentiated squamous cell carcinomas are noted for their lack of keratinization and lack of intercellular bridges. Tumors are graded with respect to the least differentiated areas. A general guideline for Fuhrman grading is given in Table-5.1 and the characteristics can be easily correlated with the images for different grades in Figure 5.1.

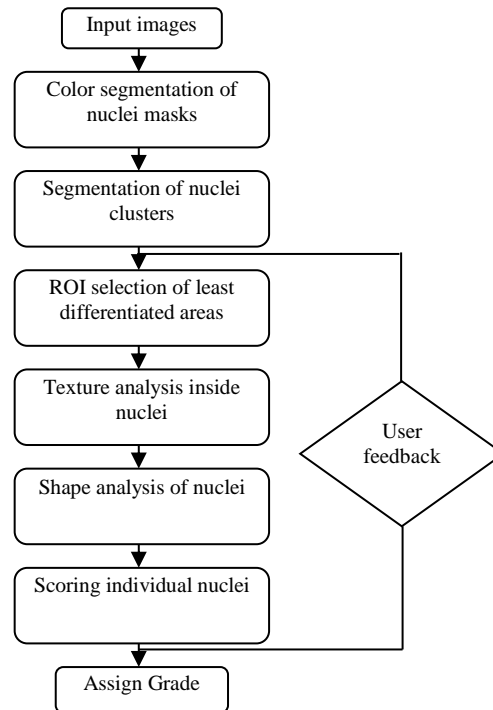
Table 5.1 General guideline for Fuhrman nuclear grading.

Grade	Nuclear Diameter in microns	Nuclear Outline	Nucleoli	Additional Features
I	10	round, regular and uniform	absent or inconspicuous	
II	15	irregular	can be seen at high power (400x), but not on low power (100x)	
III	20	irregular	prominent, easily seen at low power	
IV	20	very irregular	prominent, easily seen at low power	bizarre shapes, multi-lobed heavy chromatin clumps



**Figure 5.1** RCC images of different grades showing variation in nuclear features.

The automated grading system based on Fuhrman nuclear grade will primarily involve segmentation of the individual nuclei from clusters, determining their shape and size and analyzing the texture inside the nuclei. A proposed flowchart for the process is shown in Figure 5.2

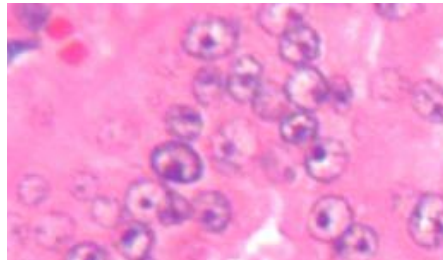


**Figure 5.2** Workflow for the proposed methodology of RCC grading

Major steps that are involved include segmenting appropriate regions (least differentiated areas), segmenting individual nuclei, texture analysis inside nuclei, shape analysis of nuclei, scoring individual nuclei and assigning combined grade to the complete image. The process accuracy may suffer due to challenges in differentiating between nuclei clusters, large malignant nuclei and artifacts. Rather than a fully automated system, we propose a decision support system involving human feedback to select regions or cells to be graded. This will improve grading accuracy for the samples with artifacts and those involving multiple grades. We also propose the grading to be carried out on a continuous scale of 1 to 4 rather than four discrete levels to utilize the precision scoring and quantification advantage available in the proposed methodology.

## **5.2 Nuclear segmentation of high grade images**

The most challenging part in an automated system for Fuhrman grading is to segment the nuclei properly otherwise the computation of nuclear characteristics can lead



to erratic results. We already have solved problem of clustered nuclei (section 3.3) but high grade nuclei have very light large regions as shown in Figure 5.3, which are unstained due to high chromatin activity. If we try to fill the regions inside nuclei using morphological techniques, neighboring nuclei being too close form clusters before complete fill is achieved. The disjoint portions of nuclei have a common characteristic generally appear to be part of same circular object. This property can be used to group these regions together to find properly segmented nuclei. Circular Hough transform (CHT) can be one of the solutions to the problem.



**Figure 5.3** Higher grade image showing lightly stained nuclei due to chromatic activity  
Nuclear Segmentation using CHT

The Hough transform can be described as a transformation of a point in the x,y plane to the parameter space. The parameter space is defined according to the shape of the object of interest.

A straight line passing through the points (x1,y1) and (x2,y2) can, in the x,y-plane, be described by:

$$y = ax+b$$

This is the equation for a straight line in the Cartesian coordinate system, where a and b represent the parameters of the line. The Hough transform for lines does not use this representation of lines, since lines perpendicular to the x-axis will have an a-value of infinity. This will force the parameter space a,b to have infinite size. Instead a line is represented by its normal which can be represented by an angle  $\theta$  and a length  $\rho$ .

$$\rho = x \cos(\theta) + y \sin(\theta)$$

The parameter space can now be spanned by  $\theta$  and  $\rho$ , where  $\theta$  will have a finite size, depending on the resolution used for  $\theta$ . The circle is actually simpler to represent in parameter space, compared to the line, since the parameters of the circle can be directly transferred to the parameter space. The equation of a circle is

$$r^2 = (x-a)^2 + (y-b)^2$$

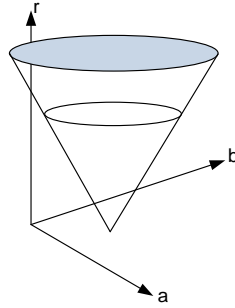
As it can be seen the circle has three parameters, r, a and b. Where a and b are the center of the circle in the x and y direction respectively and where r is the radius. The parametric representation of the circle is

$$x = a + r \cos(\theta)$$

$$y = b + r \sin(\theta)$$

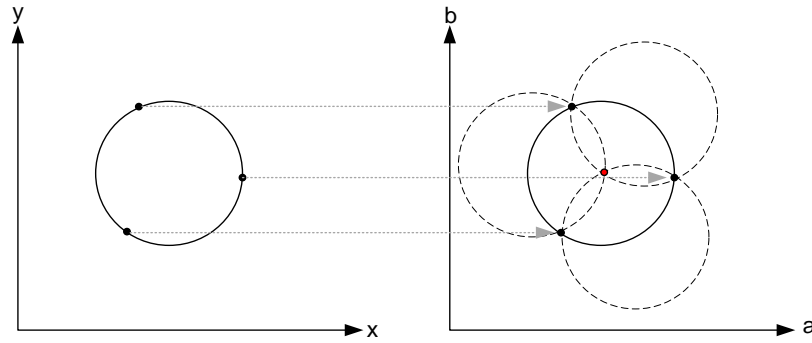
Thus the parameter space for a circle will belong to  $R^3$  whereas the line only belonged to  $R^2$ . As the number of parameters needed to describe the shape increases as well as the dimension of the parameter space R increases so does the complexity of the

Hough transform. Therefore is the Hough transform in general only considered for simple shapes with parameters belonging to  $\mathbb{R}^2$  or at most  $\mathbb{R}^3$ . In order to simplify the parametric representation of the circle, the radius can be held as a constant or limited to number of known radii.



**Figure 5.4** The parameter space used for CHT

The process of finding circles in an image using CHT is to start by finding all edges in the image. At each edge point we draw a circle with center in the point with the desired radius. This circle is drawn in the parameter space, such that our x axis is the a-value and the y axis is the b value while the z axis is the radii. At the coordinates which belong to the perimeter of the drawn circle we increment the value in our accumulator matrix which essentially has the same size as the parameter space. In this way we sweep over every edge point in the input image drawing circles with the desired radii and incrementing the values in our accumulator. When every edge point and every desired radius is used, we can turn our attention to the accumulator. The accumulator will now contain numbers corresponding to the number of circles passing through the individual coordinates. Thus the highest numbers correspond to the center of the circles in the image.

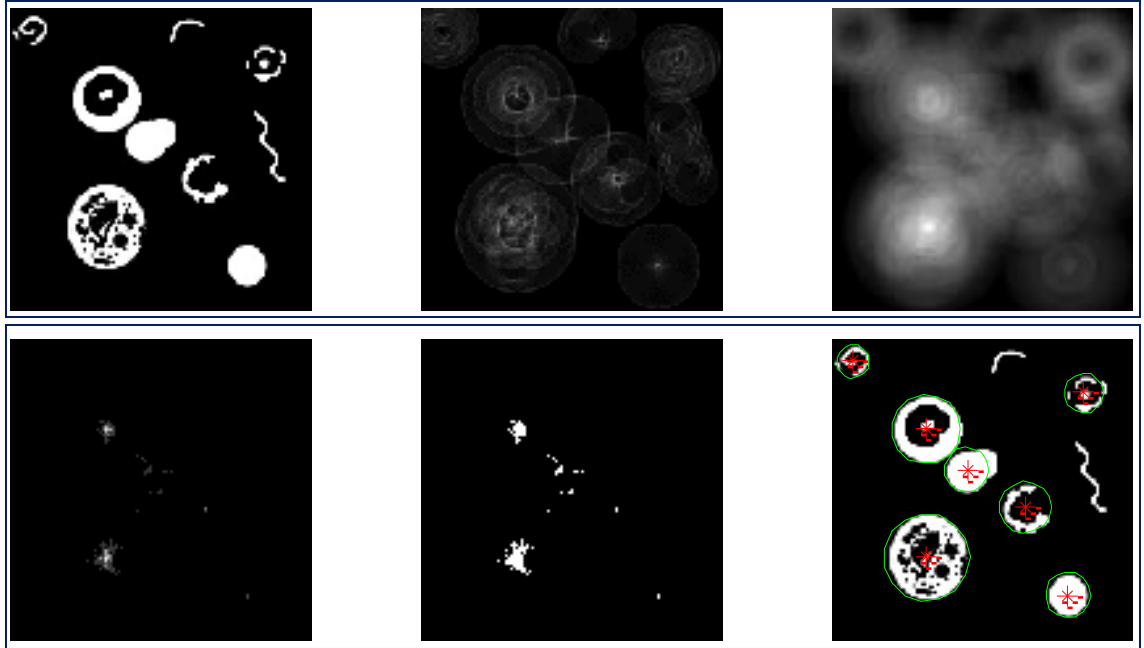


**Figure 5.5** A Circular Hough transform from the x,y-space (left) to the parameter space (right), this example is for a constant radius

We implemented this solution based on standard CHT implementation[100]. The steps in the algorithm can be listed as below:

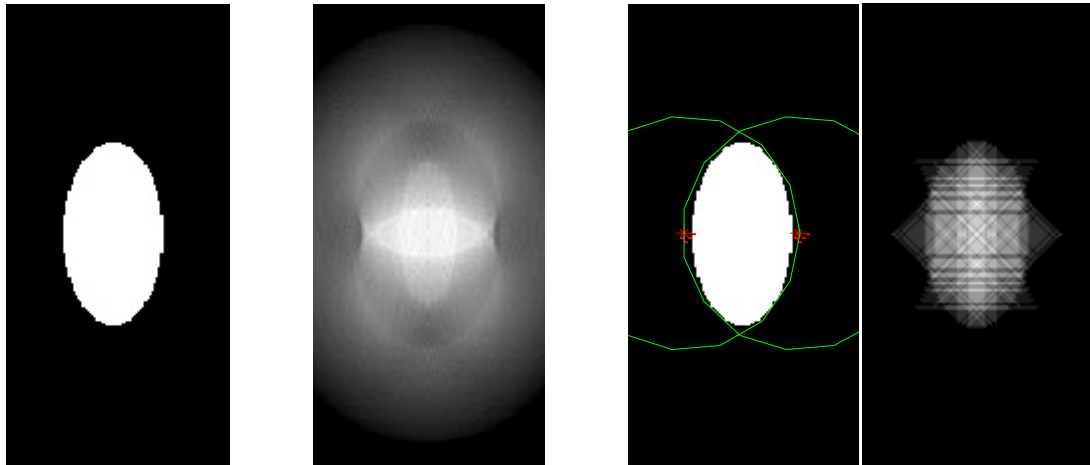
- Convert images to grayscale.
- Find edges
- Iterate for all possible radii sizes  $r$ 
  - For each edge point
  - Draw a circle with center in the edge point with radius  $r$  and increment all coordinates that the perimeter of the circle passes through in the accumulator.
  - Find one or several maxima in the accumulator
- Map the found parameters  $(r,a,b)$  corresponding to the maxima back to the original image

We started testing our implementation on a synthetic image with all possible variations. The major variation from standard implementation was only considering accumulator values above a specific threshold and using large disk (average size of nuclei) for local maxima. The algorithm segments circular objects properly and rejects unwanted objects as shown in Figure 5.6



**Figure 5.6**  $\begin{bmatrix} abc \\ def \end{bmatrix}$  (a) Synthetic image (b) Accumulator array for one specific radius (c) sum of all radii accumulator planes (d) Sum of accumulator after thresholding values (e) Binary mask of (d) (f) resultant circular models.

Another issue with CHT surfaces when the nuclei are too ecentric. In this case they can't be modeled as circular objects and may result in errors as shown in Figure 5.7(c) where a single elliptical cell has been modeled as two circular objects.



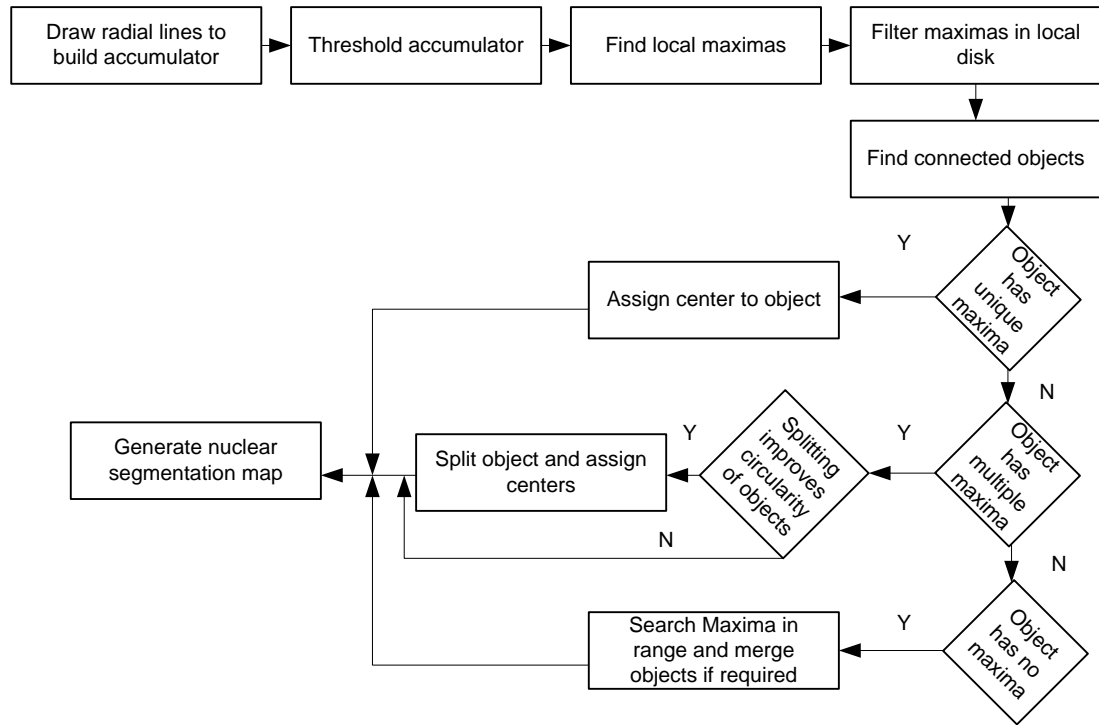
**Figure 5.7**  $[a \ b \ c \ d]$  (a) Synthetic image (b) Sum of all radii accumulator planes (c) Resultant circular models. (d) Accumulator using radial lines

Before we discuss the solution to the problem of elliptical objects we shall consider a common approach used with CHT i.e. use of gradient information. This can be a useful addition to this algorithm as it will reduce the number of points incremented in accumulator array. This can help increase processing speed and may improve accuracy.

We tested the gradient based approach and the results are similar to the one with full accumulator results with some improvement in speed.

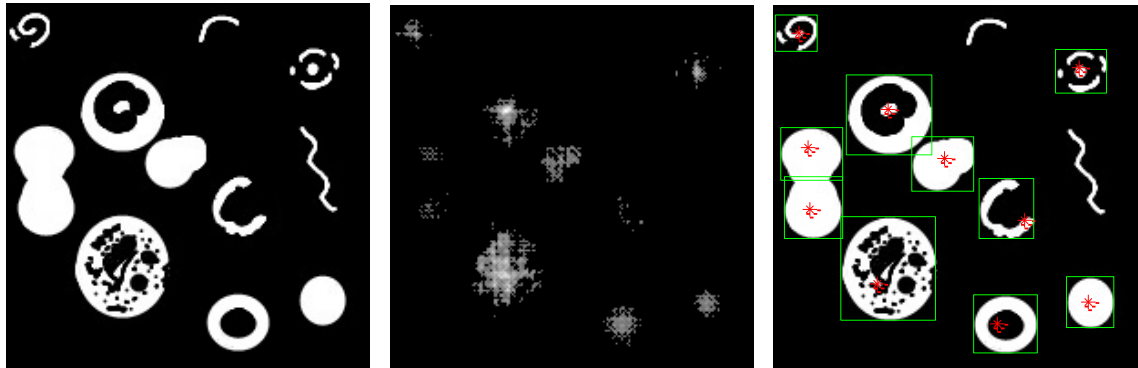
#### Nuclear segmentation using Gradient based lines.

Instead of drawing circles around each edge pixel, we can draw radial lines (to populate accumulator) in the direction of gradient as shown in Figure 5.7(d). This is more efficient approach and also takes care of non-circular objects. For real life images we need to perform objects merging and splitting to find nuclei with more accuracy. The major criterion to split the objects is their size. For the objects which aren't detected as independent nuclei, we try to find if they belong to a nucleus within their neighborhood and merge objects if needed. The overall flow chart for this algorithm is given in Figure 5.8



**Figure 5.8** Flow chart for nuclear segmentation using gradient lines based approach

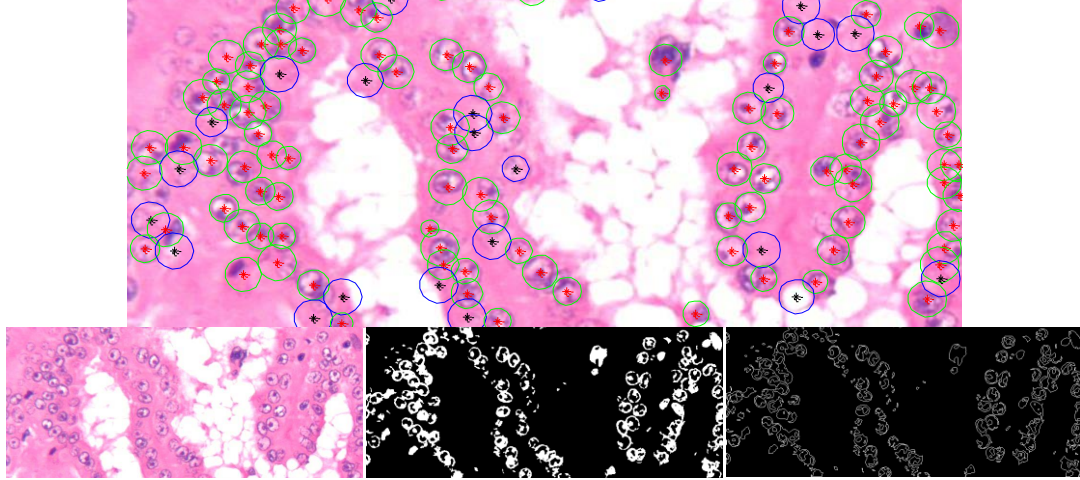
The results of this algorithm on synthetic data figure 5.9 clearly show that algorithm successfully detected nuclei and merged objects together where needed and also split a cluster into two nuclei.



**Figure 5.9** Results of gradient lines algorithm on synthetic data showing accumulator array, merging and splitting objects to detect nuclei.

Testing this scheme on one of the challenging real images in our dataset produced satisfactory results. Some regions which are not nuclei but are bounded by multiple nuclei are also detected as desired objects. The true positives are shown as green while

while false positives are shown as blue circles in Figure 5.10. These false positives can be rejected by thresholding based on percentage of stained area under the object. Figure 5.10 below shows these results.



**Figure: 5.10**  $\begin{bmatrix} a \\ bcd \end{bmatrix}$  (a) Resultant image with Green circles showing detected nuclei and blue circles showing rejected objects (b) Input image (c) Nuclear mask of input image (d) Edges used for computation of Hough transform

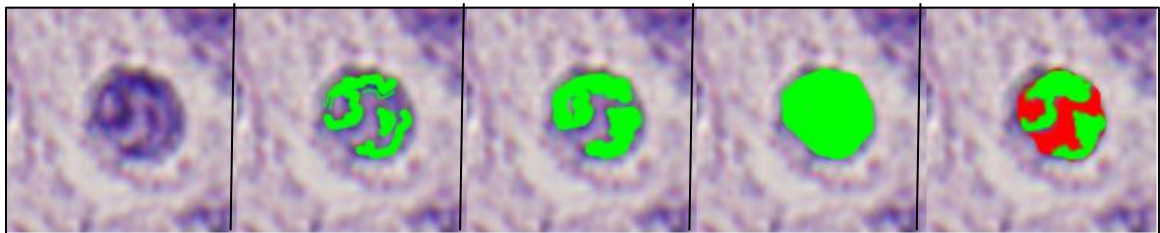
### 5.3 Feature extraction for Fuhrman grading

Fuhrman grading is primarily based on two features of the nuclei: size of the nuclei and prominence of nucleolus. The nucleolus prominence features, because of light staining, results in objects which have holes. In addition, the nuclear area segmented doesn't reflect true area or size of nuclei. We found the convex area of the segmented object to be a better measure for the nuclei size. Convex area is based on the convex hull of the object as shown in figure 5.11. The second feature related to the nucleolous is called solidity. This is an indirect measure of how solid the objects are and is computed as the proportion of the pixels in the convex hull that are also in the region. The third feature, the eccentricity, serves a dual purpose. First as a measure of circularity of the nuclei and second as a criterion to remove objects like non-split nuclei clusters and long streaks.

Taking a quantitative approach resulted in a large number of features. We selected 38 knowledge based features which can be easily co-related to the existing knowledge of pathologists. The approach helps in improved confidence in use of the quantified features for decision making. Details of these features are shown in Table 5.2:

Table 5.2 Features used for grading of RCC images

Type of features	Number of features	Most useful features
Nuclear morphology, shape and texture	12	Nuclear area, convex hull area of nuclei, filled nuclear area, unstained nuclear area, nuclear area variance, solidity, circularity, stained/unstained area ratio, area solidity ratio, edge to edge inter-nuclei distance, center to center inter-nuclei distance, and number of nuclei.
Cytoplasm features	11	
Image based features	8	
Unstained region/ background features	7	



**Figure: 5.11** [abcde](a) Nucleus image (b) Nuclear area (c) Nuclear area filled (d) Convex hull area of nuclei (e) Nuclear area unstained - red

#### 5.4 Scoring Nuclei

The pathologist, while grading ignore a lot of information like healthy cells and normal regions. They only evaluate nuclear characteristics in the least differentiated areas. Following the same guidelines, we selectively pickup nuclei for analysis. The major criterion for our selection is the maximum size of the nuclei which are expected to be seen at a given magnification and the eccentricity. If the nuclei size is too large, it means that its is a cluster which has not been split due to its shape or it may be an artifact. The



same logic applies to the eccentricity. Objects with eccentricity less than 0.5 are generally correctly identified nuclei. The second consideration based on the nuclei size is that we may ignore small sized nuclei, as pathologist do in real life, to ignore healthy and lower grade nuclei. The nuclei which are relatively large must be considered as a basis for grade prediction. At present our application is designed to select top 25% of the nuclei (based on size) for analysis.

The psuedocode for individual nuclei scoring is given below along with the explanation of variables at the end of code:

```

if  $Per_{EccN} > T_{EccPercent}$ 
    return Grade 4 // bizzare shape too many elliptical nuclei
else
    for count = 1 to  $Num_{Nuclei}$ 
        if  $N_{Ecc} < T_{Ecc}$  //confident, it is a regular nuclei
             $Flag_N(count) = set$ 
             $Score_N(count) = N_{size} * Wt_{Size} + N_{solidity} * Wt_{Solidity}$ 
        endif
    endfor
     $N_{sort} = \text{Sorted nuclei high to low}$ 
    for count = 1 to  $N_{ratioSel} * Num_{Nuclei}$ 
         $Grade = Score_{N_{sort}}(count) + Grade$ 
    end for
    return  $Grade / (N_{ratioSel} * Num_{Nuclei})$  //Avearage grade of top Nuclei
endif

```

#### **Variables:**

$T_{Ecc}$ : Threshold for rejection of nuclei based on eccentricity for individual score

$Per_{EccN}$ : Percentage of nuclei with eccentricity above  $T_{Ecc}$

$T_{\text{EccPercent}}$ : Threshold for percentage of eccentric nuclei to be considered bizarre shapes high grade image

$N_{\text{Size}}$ : Normalized nuclear size based on image resolution

$W_{\text{Size}}$ : Weightage of normalized nuclei size for score

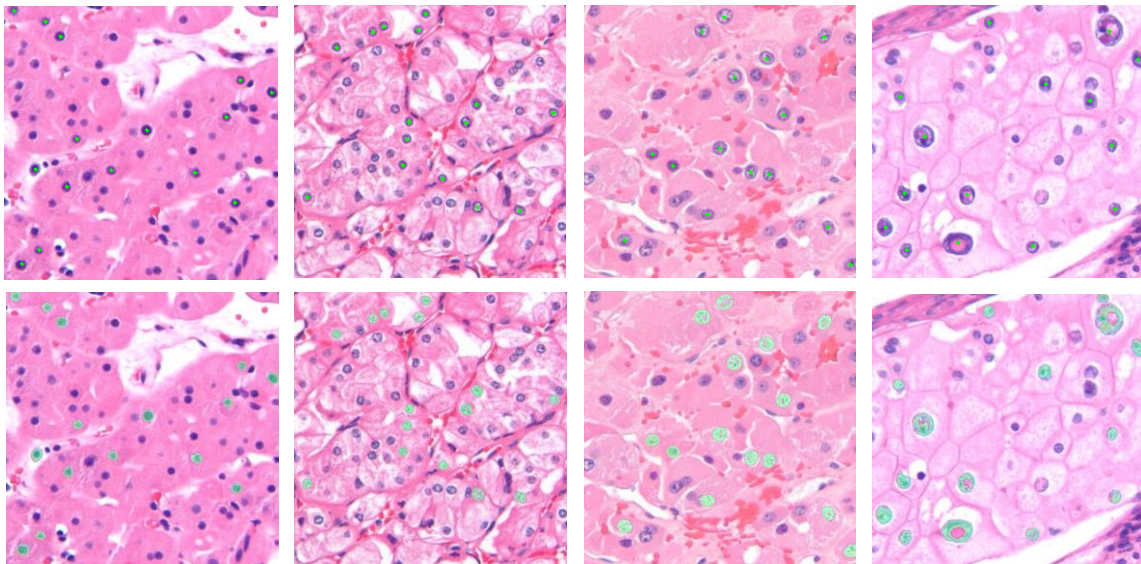
$N_{\text{Solidity}}$ : Normalized Nuclear Solidity measure of nucleolus prominence

$W_{\text{Solidity}}$ : Weightage of nuclei solidity for score

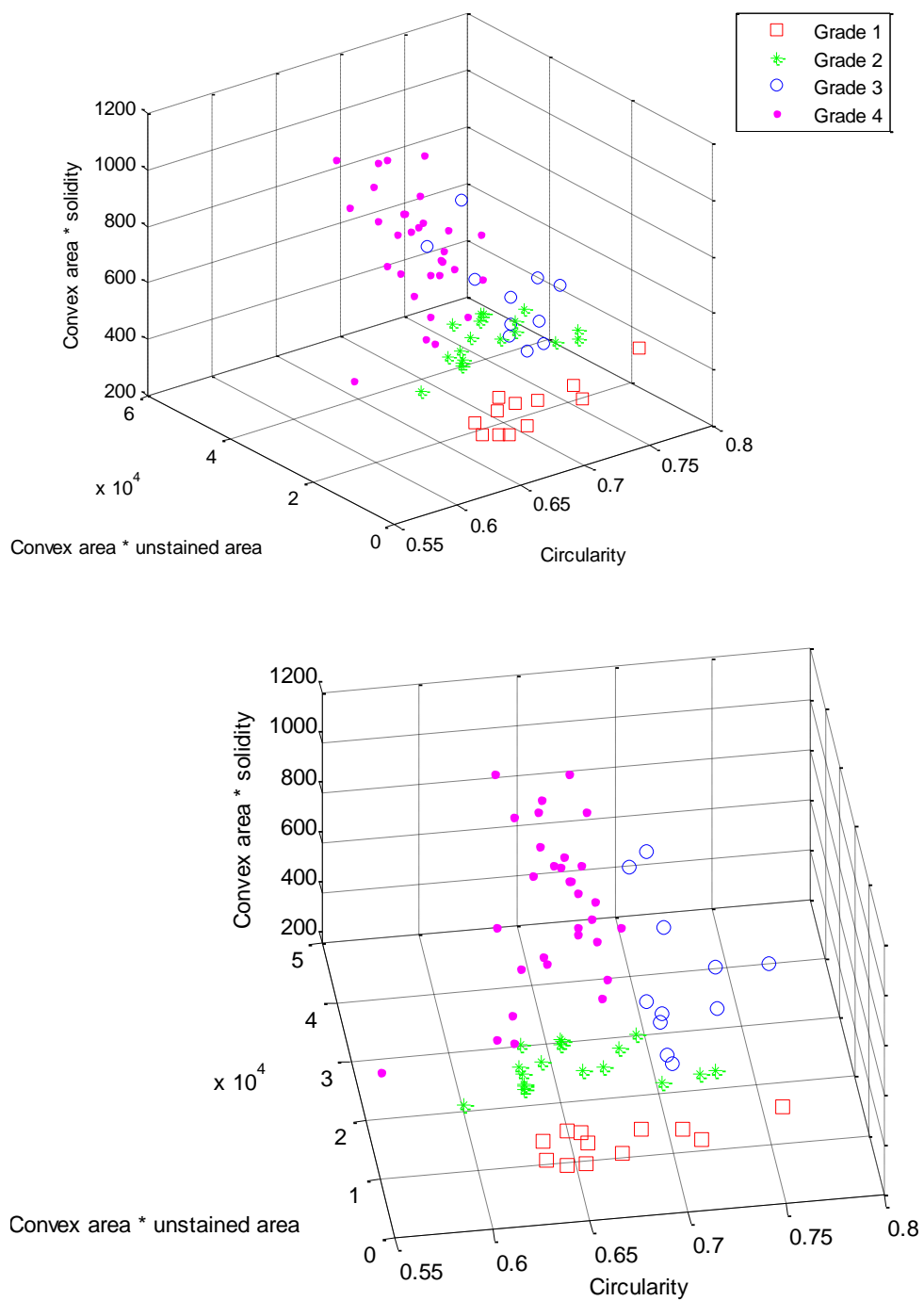
$\text{Num}_{\text{Nuclei}}$ : Number of nuclei in the image

$N_{\text{ratioSel}}$ : Percentage of top ranked nuclei considered for grading

Figure 5.12 shows the representative images of four different grades, the nuclei picked up by the algorithm for analysis and their segmentation. Analysis of selected features i.e. convex area, solidity, circularity and unstained nuclear area, shown in Figure 5.13, results in proper clustering and grade prediction adhering to fuhrman grading guidelines.



**Figure 5.12** Grade I–IV (left to right) representative images are shown with automatic selection of nuclei based on their size and eccentricity. The selected nuclei segmentation is also shown as overlay in bottom row. Holes in nuclear masks show nucleolus prominence.



**Figure 5.13** Two views of clustering of different grade images based on the selected features i.e. nuclear convex area , nuclear circularity, nuclear solidity and unstained nuclear area.

As mentioned above, after selecting top 25% nuclei based on size, we assign them individual score based on the size, eccentricity and solidity. These scores are assigned on a floating scale between 1 to 4 rather than using discrete values in conventional way. This is quite useful to compare individual nuclei and provides a usable annotation for use in the user's interface of CDSS. Use of the scoring and annotations are discussed in detail in chapter 6.

In this chapter, we reviewed the Fuhrman grading guidelines and showed that real challenge in automated grading is related to nuclear segmentation. With variations from very light staining to broken nuclear objects and overlapping nuclear cluster, the nuclear segmentation is a significant challenge. We showed how our variant of CHT successfully achieved nuclear segmentation. The algorithm also combined techniques to split and merge detected objects as detect individual nuclei. Successful detection of nuclei is followed by feature extraction and scoring of the individual nuclei leading to grading of the images. Though the knowledge extracted out of the image can lead to automated prediction of the cancer grade but our study has shown it will be more usable and more acceptable if this knowledge is used for supplementing the information available to the pathologist for use in decision making process. We designed a CDSS where the user is presented with annotated images and other information in easy to use interface to provide grading decisions in an accurate and efficient way. The details of the CDSS are discussed in the following chapter.

## **CHAPTER – VI**

### **CDSS FOR FUHRMAN GRADING**

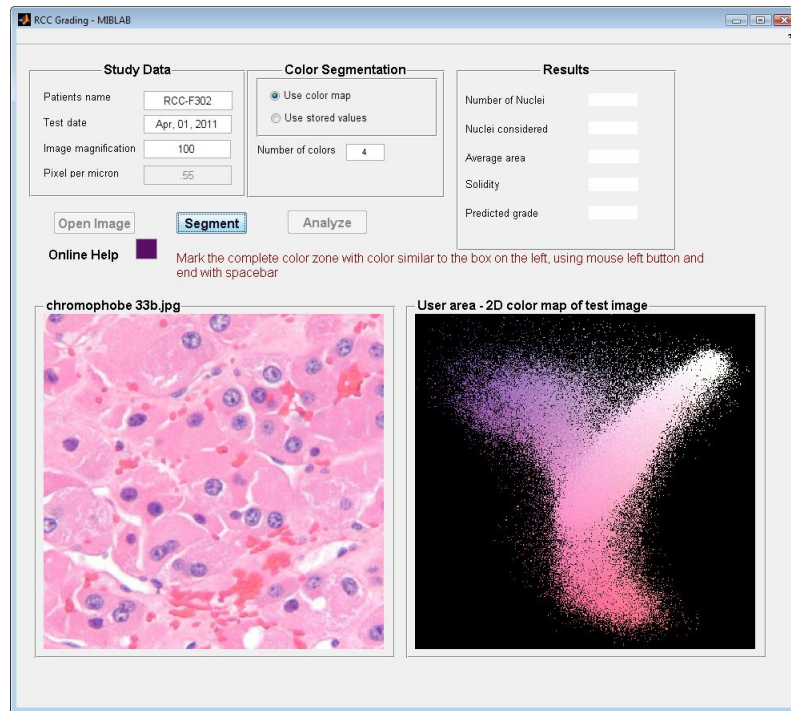
This chapter presents a Clinical Decision Support System (CDSS) for nuclear grading of Renal Cell Carcinoma (RCC). The system provides an effective way to annotate RCC tissue images with quantitative information, highlighting selected nuclei and their features, and performing automated predictions to assist pathologists in the decision making process for nuclear grading of RCC. Different image processing and analysis algorithms including color and nuclear segmentation techniques, presented in earlier chapters, can be used to reliably extract and visualize prime features which can facilitate decision making process. Testing our system on clinically challenging dataset of images, we were successfully able to extract the desired features and annotate the images for review. The system was helpful in identifying real clinical challenges and extending the clinical decision making process beyond the existing Fuhrman grading guidelines.

The investigation, analysis and interpretation of the pathological imaging data mainly depends on the pathologist's knowledge, experience and his subjective view about the data. As Computer Aided Diagnostic (CAD) tools can help reduce this subjectivity and the inter-user and intra-user variability, they are gaining some acceptance among clinicians [3, 4]. The clinicians, however, prefer systems which are flexible and take into account their individuality when necessary by providing some control rather than fully automated system [101]. Studies[34, 35] have shown improvement in practitioner's performance and patient's outcome for CDSS which account for practitioner's perspective and integrate it in the workflow. Therefore, to be able to introduce CDSS in health care, we need to understand users' perspectives and preferences on the new information technology. This forms as the basis for this CDSS where we target to present the quantitative information acquired through the image

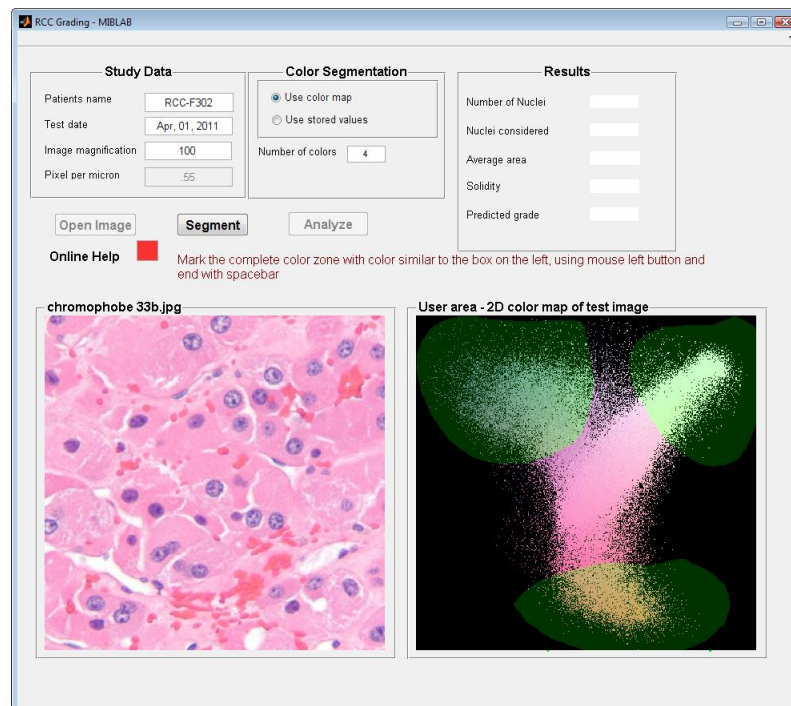
analysis, annotate the images and provide suitable visualization which can facilitate the process of decision making in a clinical setting.

### **6.1 Graphical user interface (GUI)**

We designed a GUI which combines the power of 2D color map based segmentation (section 3.1.2), gradient lines based nuclear segmentation (section 5.1) and the feature extraction for Fuhrman grading (section 5. 2) to provide an interface to the user to interact and test the images. The GUI allows users to open all common types of images including jpeg, tiff etc. The user is presented with a 2D color map with step by step instructions (figure 6.1). The user marks zones using mouse pointer which are perceived to belong to a single color class (figure 6.2). The number of color classes can be selected using color segmentation data panel. The color segmentation algorithm runs in the background and color segmentation results are presented to the user (figure 6.3). Correct segmentation enhances user's confidence to proceed to the analysis step. The analysis step computes features and generates vital statistics which are shown in the results panel. Top ranked nuclei, based on the grading score, are highlighted using green overlay over the input image. The predicted grade for the overall image is shown in results panel (figure 6.4). The final decision made by the user is recorded in the feedback panel (figure 6.4) along with the system performance evaluation and specific comments related to the image under study.

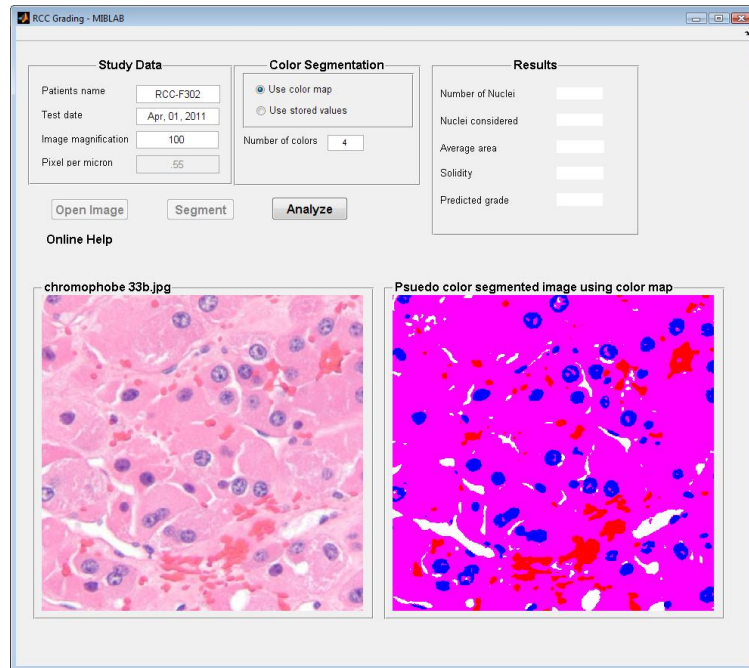


**Figure 6.1** GUI showing image loading, entry for basic parameters, 2D segmentation map and user instructions.



**Figure 6.2** GUI showing zone marking by user for segmenation





**Figure 6.3** GUI showing segmenation results in psuedo color.



**Figure 6.4** GUI showing top ranked nuclei(overlaid in green) based on their grading score(black). The holes in the segmentation masks show nucleolus prominence. Major statistics along with the predicted grade are shown in results panel. Feedback panel records final decision, system performance and specific comments.



## 6.2 Time analysis

The GUI shown in section 6.1 has two major components i.e. the color segmentation and the feature extraction and analysis component. The interactive use of color segmentation component is not required every time and is only recommended once for a single batch of images stained and captured under identical conditions. Moreover, for whole slide scans which may be captured using slide scanners or by using image mosaicking to join multiple images, a single representative image can be used to capture user color perception. The time taken by this component constitutes of time to generate the color map (less than 1 sec), the user interaction time (less than 40 secs.) and generating segmentation results(less than 1 sec) is dominated by the user interaction part. Detailed time analysis results for this component are shown in the user study (section 3.1.2, table 3.2 and figure 3.10).

The feature extraction and analysis component which works on the color segmented images is designed to work on image tiles. Smaller images e.g. 1000x1000 pixels can be processed as a single image or single tile. The whole slide images are split in multiple tiles (sub-images) and each tile is analyzed independently. The statistics from these tiles are merged to generate combined image statistics. The time taken by large images is directly proportional to the ratio of its area to the tile area. Table 6.1 shows time analysis of some example images.

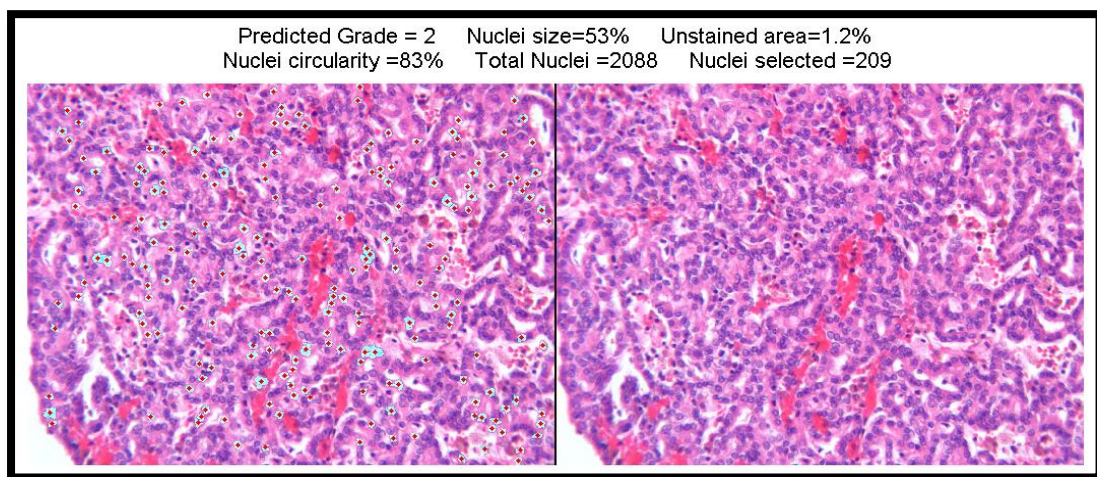
Table 6.1 Time analysis showing time required for different components for image examples of different sizes.

<b>Image Size</b>	<b>Tile size</b>	<b>Number of tiles</b>	<b>User interaction time other than segmentation time</b>	<b>Color segmentation time</b>	<b>Feature extraction and analysis time</b>	<b>Total time (secs)</b>
600x600	600x600	1	10.5secs	21.32 secs	6.88 secs	38.7 secs
1600x1200	800x600	4	12.3 secs	25.32 secs	30.5 secs	68.12 secs

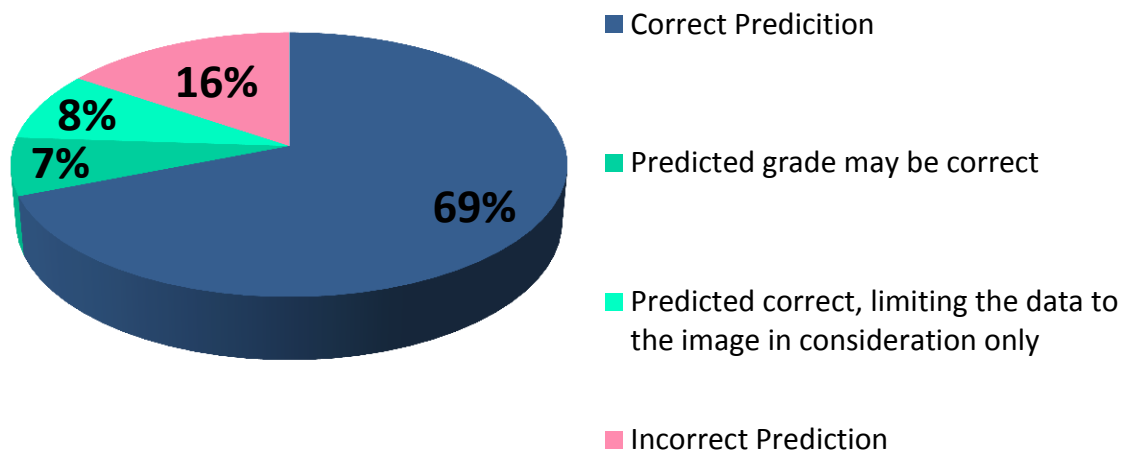
### 6.3 Grade predictions and decisions

At present the system is targeted towards a decision support system. The system can be used as a second more objective reader to aid the pathologist in making final diagnostic decision. For system evaluation, we used a data set of 72 images which were initially graded by a urologic pathologist. However, the initial grading was not done entirely based on the view covered in the image. The urologic pathologist had whole slide available for analysis and other regions of the tissue were considered when a single view wasn't good enough to base their decision. The second limitation of the system was that it was based on the Fuhrman grading guideline which is inherently limited to basic knowledge, while the pathologists use their experience and knowledge for decision making which is way beyond the stated guidelines like Fuhrman grading. We generated annotations for these images (Figure 6.5), marked nuclei which were selected for grade statistics and assigned the new grade labels based on quantitative analysis. These annotated images were again reviewed by the urologic pathologist. During review, the results (Figure 6.6) show that 69% of the urologic pathologist's initial assigned grades were correctly predicted by the system, 7% of images lead to re-evaluation of grade with

a potential to change the original prediction with system recommendation. 8% of images can be considered as correct prediction if the data is restricted to the information available in the image; however the initial predictions by the urologic pathologist were based on other regions of tissue which was not available to the system. 16% of the images were considered as wrong predictions, primarily due to the system design being restricted to the implementation of Fuhrman grading guidelines. Our average prediction error on a four tier grade system was close to 0.6 and 98.6% of the predictions were less than 1 grade level from the pathologist's assigned labels.



**Figure 6.5** Original and annotated images are shown which can be used for pathologists review



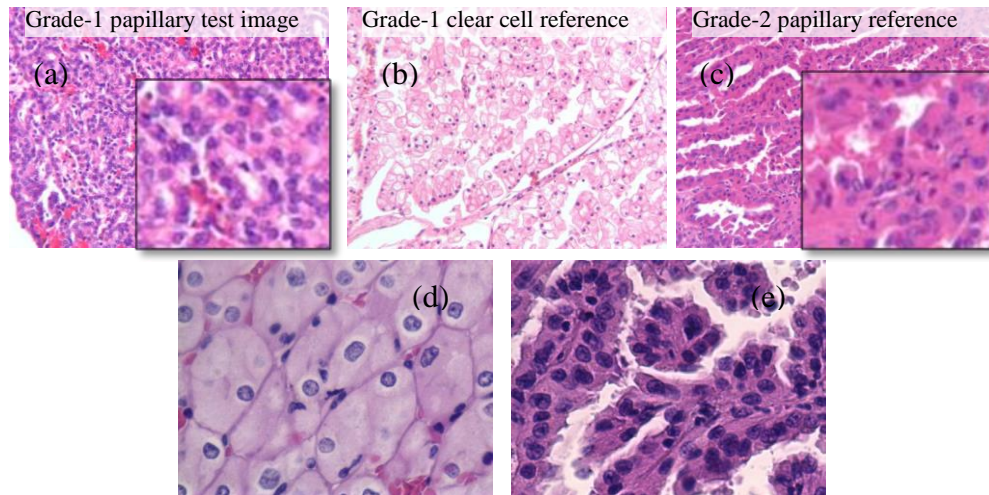
**Figure 6.6** Prediction results for the case study of 72 RCC images

#### 6.4 Feedback and results

The feedback collected during the review of the annotated images revealed certain limitations of the Fuhrman grading system which were not specifically mentioned in the guideline but have become part of the decision making process for the expert pathologists. A few of these limitations are.

- a. Fuhrman grading was primarily designed for the clear cell carcinoma being the most frequently occurring histopathologic subtype of RCC. Different subtypes have unique characteristics and subtype consideration is an important aspect which is not taken care off in Fuhrman grading. Figure 6.7 shows that in case of papillary carcinoma the nuclear size is generally larger than the one in other subtypes.
- b. Pleomorphism, which is measure of nuclear variance, is another important consideration not provided in Fuhrman guideline. For nuclei with similar average size and similar nucleolus prominence the tissue exhibiting pleomorphism will tend to have higher grade as shown in figure 6.7

- c. Another important feature to determine the grade is bi-nucleation which is a measure of presence of multiple nuclei in a single cell. H&E stain, generally used for RCC tissue images, isn't targeting cell membranes. It is very hard to see cell boundaries and therefore determine whether closely located nuclei are part of single cell or independent cells. Figure 6.7(d&e) shows examples with and without prominent cell membranes. In case it can be established that bi-nucleation is present, the image is categorized in higher grade.



**Figure 6.7** (a&b) Nuclei size for papillary carcinoma is larger than the one in corresponding clear cell reference. (a&c) For same size and nucleolus prominence image exhibiting pleomorphism shows higher grade. (d) Cell structure is clear because of prominent cell boundaries. (e) Nuclei are clustered together but cell boundaries aren't visible and bi-nucleation cannot be determined.

## **Summary**

Our case study of RCC grading has shown that automated analysis where very high prediction accuracies are hard to achieve due to complexity and variability of the clinical data, the CDSS solutions may be a better approach rather than conventional approach of using fully automated CAD tools. The clinicians can still benefit from useful semantic knowledge; enhanced visualizations and predictions in a flexible and usable interface. It makes the decision making process easier and adds to confidence in clinical decisions. During the developmental phase, the instances where the system predictions contradict the clinician's observations can prompt detailed investigation of more data and may improve decisions after review. Clinicians recorded feedback and annotated image database can be a vital resource to improve the knowledge base and subsequent predictions thereby leading to a better and mature CDSS which can be fully integrated in the clinical workflow.

## **CHAPTER – VII**

### **CONCLUSION**

The successful case study of RCC subtype classification and grading has shown promise for our CAD tools to be used as a more general pathology system especially in clinical setting as CDSS. With the ability to process differently stained images, ability to segment images efficiently and accurately and the flexibility to customize features already built-in, the system can readily be used for the classification and grading of other cancer imaging data.

In order to ensure our tool's universal clinical use, more work needs to be done. In addition to the features already discussed here, a whole array of features may be added to the system including gabor filters, phase congruence analysis and fractal vectors. All these features will be available for the user to select from depending on his or her preference and the perceived importance during the diagnosis. Aside from adding new feature extraction, nuclear and cellular segmentation is open to exploration. New techniques for merging and splitting cell clusters can help better segmentation of the objects which will in turn result in better feature statistics and improved classification and grading results.

It should be kept in mind, however that a clinical implementation would require much more training and validation. In particular, additional effort is required to cater to the variation among image acquisition systems, tissue collection, and staining protocols.

Although computer aided diagnostic tools can never replace the expertise of a trained pathologist, they can definitely assist the pathologists and improve upon the status quo by increasing accuracy and reducing subjectivity. Furthermore, the power of pattern recognition can also be leveraged to find and extract features from these tissue images that are beyond human visual perception.

It is critical to note that the paradigm of cancer diagnostics based mainly on histopathology will soon have to evolve. True[89] and Gao[20] have contended that since histological patterns of cancer are not directly correlated with the underlying molecular profile that is responsible for cancer progression, new optical image technologies like Quantum dots should be used to provide further insight. They state that “With new molecular profiling technologies, it should be possible to read the molecular signatures of an individual patient’s tumor and correlate a panel of tissue biomarkers with clinical outcome and personalized therapy”. Our system can readily analyze and quantify such images and their properties and can integrate that information with the underlying histopathological information. The advent of these molecular profiling techniques coupled with pathological CDSS promises the ability to correlate the histological patterns with the specific biomarker profile thereby leading to better clinical diagnosis and represents a new horizon for the field of molecular pathology.



## REFERENCES

- [1] A. Jemal, R. Siegel, E. Ward, T. Murray, J. Xu, and M. J. Thun, "Cancer statistics, 2007," *CA: a cancer journal for clinicians*, vol. 57, no. 1, pp. 43-66, 2007.
- [2] A. C. Society. "Cancer Facts and Figures 2010. Atlanta, Ga: American Cancer Society," Mar 2011, 2011; <http://www.cancer.org/acs/groups/content/@nho/documents/document/acspc-024113.pdf>.
- [3] E. R. Farmer, R. Gonin, and M. P. Hanna, "Discordance in the histopathologic diagnosis of melanoma and melanocytic nevi between expert pathologists," *Human pathology*, vol. 27, no. 6, pp. 528-531, 1996.
- [4] A. B. Ackerman, "Discordance among expert pathologists in diagnosis of melanocytic neoplasms," *Human pathology*, vol. 27, no. 11, pp. 1115-1116, 1996.
- [5] K. Fleming, "Evidence-based pathology," *The Journal of Pathology*, vol. 179, no. 2, pp. 127, 1996.
- [6] J. P. Thiran, and B. Macq, "Morphological feature extraction for the classification of digital images of cancerous tissues," *Biomedical Engineering, IEEE Transactions on*, vol. 43, no. 10, pp. 1011-1020, 2002.
- [7] R. F. Walker, P. Jackway, B. Lovell, and I. Longstaff, "Classification of cervical cell nuclei using morphological segmentation and textural feature extraction." pp. 297-301.
- [8] Y. Jiang, R. Nishikawa, R. Schmidt, C. Metz, M. Giger, and K. Doi, "Improving breast cancer diagnosis with computer-aided diagnosis," *Academic Radiology*, vol. 6, no. 1, pp. 22-33, 1999.
- [9] M. Roula, J. Diamond, A. Bouridane, P. Miller, and A. Amira, "A multispectral computer vision system for automatic grading of prostatic neoplasia." pp. 193-196.
- [10] J. Diamond, N. Anderson, P. Bartels, R. Montironi, and P. Hamilton, "The use of morphological characteristics and texture analysis in the identification of tissue composition in prostatic neoplasia," *Human pathology*, vol. 35, no. 9, pp. 1121-1131, 2004.
- [11] N. Esgiar, R. Naguib, B. Sharif, M. Bennett, and A. Murray, "Microscopic image analysis for quantitative measurement and feature identification of normal and cancerous colonic mucosa," *Information Technology in Biomedicine, IEEE Transactions on*, vol. 2, no. 3, pp. 197-203, 2002.
- [12] A. Esgiar, R. Naguib, B. Sharif, M. Bennett, and A. Murray, "Fractal analysis in the detection of colonic cancer images," *Information Technology in Biomedicine, IEEE Transactions on*, vol. 6, no. 1, pp. 54-58, 2002.
- [13] P. Hamilton, P. Bartels, D. Thompson, N. Anderson, R. Montironi, and J. Sloan, "Automated location of dysplastic fields in colorectal histology using image texture analysis," *The Journal of Pathology*, vol. 182, no. 1, pp. 68-75, 1997.
- [14] M. Amin, P. Tamboli, J. Javidan, H. Stricker, M. Venturina, A. Deshpande, and M. Menon, "Prognostic impact of histologic subtyping of adult renal epithelial neoplasms," *Am J Surg Pathol*, vol. 26, no. 3, pp. 281-291, 2002.

- [15] J. C. Cheville, C. M. Lohse, H. Zincke, A. L. Weaver, and M. L. Blute, "Comparisons of outcome and prognostic features among histologic subtypes of renal cell carcinoma," *The American Journal of Surgical Pathology*, vol. 27, no. 5, pp. 612, 2003.
- [16] Wikipedia. "Staining," 2011; <http://en.wikipedia.org/wiki/Staining>.
- [17] M. Z. Bruckner. "Basic Cellular Staining," March 06, 2011 2011; [http://serc.carleton.edu/microbelife/research\\_methods/microscopy/cellstain.html](http://serc.carleton.edu/microbelife/research_methods/microscopy/cellstain.html)
- [18] M. Bruchez, M. Moronne, P. Gin, S. Weiss, and A. P. Alivisatos, "Semiconductor nanocrystals as fluorescent biological labels," *Science*, vol. 281, no. 5385, pp. 2013, 1998.
- [19] W. C. W. Chan, and S. Nie, "Quantum dot bioconjugates for ultrasensitive nonisotopic detection," *Science*, vol. 281, no. 5385, pp. 2016, 1998.
- [20] X. Gao, and S. Nie, "Molecular profiling of single cells and tissue specimens with quantum dots," *TRENDS in Biotechnology*, vol. 21, no. 9, pp. 371-373, 2003.
- [21] R. Nisman, G. Dellaire, Y. Ren, R. Li, and D. P. Bazett-Jones, "Application of quantum dots as probes for correlative fluorescence, conventional, and energy-filtered transmission electron microscopy," *Journal of Histochemistry and Cytochemistry*, vol. 52, no. 1, pp. 13, 2004.
- [22] Y. Xing, Q. Chaudry, C. Shen, K. Kong, H. Zhau, L. Chung, J. Petros, R. O'Regan, M. Yezhelyev, and J. Simons, "Bioconjugated quantum dots for multiplexed and quantitative immunohistochemistry," *Nature Protocols*, vol. 2, no. 5, pp. 1152-1165, 2007.
- [23] E. S. Berner, *Clinical Decision Support Systems, Theory and Practice*: Springer, 2006.
- [24] D. F. Sittig, A. Wright, J. A. Osherooff, B. Middleton, J. M. Teich, J. S. Ash, E. Campbell, and D. W. Bates, "Grand Challenges in Clinical Decision Support v10," *Journal of biomedical informatics*, vol. 41, no. 2, pp. 387, 2008.
- [25] A. Depeursinge, D. Racocanu, J. Iavindrasana, G. Cohen, A. Platon, P. A. Poletti, and H. Müller, "Fusing visual and clinical information for lung tissue classification in high-resolution computed tomography," *Artificial intelligence in medicine*, vol. 50, no. 1, pp. 13-21, 2010.
- [26] H. J. Lee, S. I. Hwang, S. Han, S. H. Park, S. H. Kim, J. Y. Cho, C. G. Seong, and G. Choe, "Image-based clinical decision support for transrectal ultrasound in the diagnosis of prostate cancer: comparison of multiple logistic regression, artificial neural network, and support vector machine," *European radiology*, vol. 20, no. 6, pp. 1476-1484, 2010.
- [27] D. Demner-Fushman, S. Antani, and G. R. Thoma, "Automatically finding images for clinical decision support," *Data Mining Workshops, 2007. ICDM Workshops 2007. Seventh IEEE International Conference on*, pp. 139-144, 2007.
- [28] D. Comaniciu, P. Meer, and D. J. Foran, "Image-guided decision support system for pathology," *Machine Vision and Applications*, vol. 11, no. 4, pp. 213-224, 1999.
- [29] K. Congdon, "Top 10 IT Trends for 2012," *Health care technology online* Feb 2012.
- [30] "Clinical Decision Support Systems," <http://healthinformatics.wikispaces.com/Clinical+Decision+Support+Systems>.

- [31] W. Chen, P. Meer, B. Georgescu, W. He, L. A. Goodell, and D. J. Foran, "Image mining for investigative pathology using optimized feature extraction and data fusion," *Computer methods and programs in biomedicine*, vol. 79, no. 1, pp. 59-72, 2005.
- [32] F. Chiarugi, S. Colantonio, D. Emmanouilidou, D. Moroni, F. Perticone, A. Sciacqua, and O. Salvetti, "ECG and echocardiography processing for decision support in heart failure," *Computers in Cardiology*, 2008, pp. 649-652, 2008.
- [33] O. Sertel, J. Kong, H. Shimada, U. Catalyurek, J. H. Saltz, and M. N. Gurcan, "Computer-aided prognosis of neuroblastoma on whole-slide images: Classification of stromal development," *Pattern Recognition*, vol. 42, no. 6, pp. 1093-1103, 2009.
- [34] A. X. Garg, N. K. J. Adhikari, H. McDonald, M. P. Rosas-Arellano, P. Devereaux, J. Beyene, J. Sam, and R. B. Haynes, "Effects of computerized clinical decision support systems on practitioner performance and patient outcomes," *JAMA: the journal of the American Medical Association*, vol. 293, no. 10, pp. 1223-1238, 2005.
- [35] M. W. M. Jaspers, M. Smeulers, H. Vermeulen, and L. W. Peute, "Effects of clinical decision-support systems on practitioner performance and patient outcomes: a synthesis of high-quality systematic review findings," *Journal of the American Medical Informatics Association*, vol. 18, no. 3, pp. 327-334, 2011.
- [36] H. Varonen, T. Kortteisto, and M. Kaila, "What may help or hinder the implementation of computerized decision support systems (CDSSs): a focus group study with physicians," *Family practice*, vol. 25, no. 3, pp. 162-167, 2008.
- [37] B. Kaplan, "Evaluating informatics applications—clinical decision support systems literature review," *International journal of medical informatics*, vol. 64, no. 1, pp. 15-37, 2001.
- [38] G. Kong, D. L. Xu, and J. B. Yang, "Clinical decision support systems: a review on knowledge representation and inference under uncertainties," *International Journal of Computational Intelligence Systems*, vol. 1, no. 2, pp. 159-167, 2008.
- [39] L. Liotta, and E. Petricoin, "Molecular profiling of human cancer," *Nature Reviews Genetics*, vol. 1, no. 1, pp. 48-56, 2000.
- [40] D. Cross, and J. K. Burmester, "The promise of molecular profiling for cancer identification and treatment," *Clinical Medicine & Research*, vol. 2, no. 3, pp. 147, 2004.
- [41] J. Ioannidis, "Is molecular profiling ready for use in clinical decision making?," *The oncologist*, vol. 12, no. 3, pp. 301, 2007.
- [42] Q. Chaudry, K. Kong, T. Ahearn, V. Cohen, R. Bostick, and M. Wang, "An integrated image quantification system for colorectal cancer risk assessment using quantum dots and molecular profiling." pp. 1280-1283.
- [43] R. M. Bostick, K. Y. Kong, T. U. Ahearn, Q. Chaudry, V. Cohen, and M. D. Wang, "Detecting and quantifying biomarkers of risk for colorectal cancer using quantum dots and novel image analysis algorithms." pp. 3313-3316.
- [44] K. Wu, D. Gauthier, and M. Levine, "Live cell image segmentation," *Biomedical Engineering, IEEE Transactions on*, vol. 42, no. 1, pp. 1-12, 2002.
- [45] S. Umbaugh, R. Moss, W. Stoecker, and G. Hance, "Automatic color segmentation algorithms-with application to skin tumor feature identification,"

- IEEE Engineering in Medicine and Biology Magazine*, vol. 12, no. 3, pp. 75-82, 1993.
- [46] P. Gejgus, and J. Placek, "Skin color segmentation method based on mixture of gaussians and its application in learning system for finger alphabet." pp. 1-6.
  - [47] Y. Fang, and T. Tan, "A novel adaptive colour segmentation algorithm and its application to skin detection." pp. 23-31.
  - [48] D. Nikolaev, and P. Nikolayev, "Linear color segmentation and its implementation," *Computer Vision and Image Understanding*, vol. 94, no. 1-3, pp. 115-139, 2004.
  - [49] S. Waheed, R. Moffitt, Q. Chaudry, A. Young, and M. Wang, "Computer Aided Histopathological Classification of Cancer Subtypes." pp. 503-508.
  - [50] K. Fu, and J. Mui, "A survey on image segmentation," *Pattern Recognition*, vol. 13, no. 1, pp. 3-16, 1981.
  - [51] R. Haralick, and L. Shapiro, "Image segmentation techniques," *Computer vision, graphics, and image processing*, vol. 29, no. 1, pp. 100-132, 1985.
  - [52] N. Pal, and S. Pal, "A review on image segmentation techniques," *Pattern Recognition*, vol. 26, no. 9, pp. 1277-1294, 1993.
  - [53] A. Weeks, and G. Hague, "Color segmentation in the HSI color space using the K-means algorithm." p. 143.
  - [54] Q. Chaudry, S. Raza, A. Young, and M. Wang, "Automated Renal Cell Carcinoma Subtype Classification Using Morphological, Textural and Wavelets Based Features," *Journal of Signal Processing Systems*, vol. 55, no. 1, pp. 15-23, 2009.
  - [55] P. Schmid, and S. Fischer, "Colour segmentation for the analysis of pigmented skin lesions." pp. 688-692.
  - [56] C. Healey, and J. Enns, "A perceptual colour segmentation algorithm."
  - [57] E. Vazquez, R. Baldrich, J. Vazquez, and M. Vanrell, "Topological histogram reduction towards colour segmentation," *Pattern Recognition and Image Analysis*, pp. 55-62, 2007.
  - [58] J. Yi, J. Park, J. Kim, and J. Choi, "Robust skin color segmentation using a 2D plane of RGB color space," *Computer and Information Sciences-ISCIS 2003*, pp. 413-420, 2003.
  - [59] J. Gao, A. Kosaka, and A. Kak, "Interactive color image segmentation editor driven by active contour model." pp. 245-249.
  - [60] A. Kapelner, P. Lee, and S. Holmes, "An interactive statistical image segmentation and visualization system." pp. 81-86.
  - [61] P. Holting, and C. Wählby, "Easy-to-use object selection by color space projections and watershed segmentation," *Image Analysis and Processing-ICIAP 2005*, pp. 269-276, 2005.
  - [62] B. Maxwell, "A Physics-Based Approach to Interactive Segmentation." pp. 656-656.
  - [63] P. Loo, and C. Tan, "Adaptive region growing color segmentation for text using irregular pyramid," *Document Analysis Systems VI*, pp. 103-106, 2004.
  - [64] SharkD. "HSV color solid cylinder.png," [http://commons.wikimedia.org/wiki/File:HSV\\_color\\_solid\\_cylinder.png](http://commons.wikimedia.org/wiki/File:HSV_color_solid_cylinder.png).

- [65] D. MacKay, *Information Theory, Pattern Recognition and Neural Networks*: Cambridge University Press, 2003.
- [66] S. Osher, and N. Paragios, *Geometric level set methods in imaging, vision, and graphics*: Springer-Verlag New York Inc, 2003.
- [67] M. Haindl, and S. Mikes. "Texture segmentation benchmark," <http://mosaic.utia.cas.cz/>.
- [68] Mathworks. "Specifying a Region of Interest," 1/11/2011, 2011; <http://www.mathworks.com/help/toolbox/images/f19-13234.html>.
- [69] E. Mortensen, and W. Barrett, "Intelligent scissors for image composition." pp. 191-198.
- [70] W. Barrett, L. Reese, and E. Mortensen, "Intelligent segmentation tools." pp. 217-220.
- [71] H. Refai, L. Li, T. Teague, and R. Naukam, "Automatic count of hepatocytes in microscopic images." pp. 1101-1104.
- [72] C. Di Rubeto, A. Dempster, S. Khan, and B. Jarra, "Segmentation of blood images using morphological operators." pp. 397-400.
- [73] E. Glory, V. Meas-Yedid, G. Stamon, C. Pinset, and J. Olivo-Marin, "Automated image-based screening of cell cultures for cell therapy." pp. 259-262.
- [74] O. Schmitt, and M. Hasse, "Radial symmetries based decomposition of cell clusters in binary and gray level images," *Pattern Recognition*, vol. 41, no. 6, pp. 1905-1923, 2008.
- [75] B. Nilsson, and A. Heyden, "Segmentation of dense leukocyte clusters." pp. 221-227.
- [76] P. Thevenaz, and M. Unser, "Snakuscles," *Image Processing, IEEE Transactions on*, vol. 17, no. 4, pp. 585-593, 2008.
- [77] W. Weixing, and S. Hao, "Cell Cluster Image Segmentation on Form Analysis." pp. 833-836.
- [78] X. Bai, C. Sun, and F. Zhou, "Touching Cells Splitting by Using Concave Points and Ellipse Fitting." pp. 271-278.
- [79] Q. Chaudry, S. Raza, Y. Sharma, A. Young, and M. Wang, "Improving renal cell carcinoma classification by automatic region of interest selection." pp. 1-6.
- [80] P. Soille, *Morphological image analysis: principles and applications*: Springer-Verlag New York, Inc. Secaucus, NJ, USA, 2003.
- [81] A. Fitzgibbon, M. Pilu, and R. Fisher, "Direct least square fitting of ellipses," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 21, no. 5, pp. 476-480, 2002.
- [82] Wikipedia. "Feature Extraction," [http://en.wikipedia.org/wiki/Feature\\_extraction](http://en.wikipedia.org/wiki/Feature_extraction).
- [83] Wikipedia. "Feature Selection," [http://en.wikipedia.org/wiki/Feature\\_selection](http://en.wikipedia.org/wiki/Feature_selection).
- [84] Wikipedia. "Statistical classification," [http://en.wikipedia.org/wiki/Classification\\_%28machine\\_learning%29](http://en.wikipedia.org/wiki/Classification_%28machine_learning%29).
- [85] S. Kothari, Q. Chaudry, and M. Wang, "Extraction of informative cell features by segmentation of densely clustered tissue images." pp. 6706-6709.
- [86] M. Hauta-Kasari, J. Parkkinen, T. Jaaskelainen, and R. Lenz, "Generalized co-occurrence matrix for multispectral texture analysis." pp. 785-789.
- [87] P. Bamford, and B. Lovell, "A water immersion algorithm for cytological image segmentation." pp. 75-79.

- [88] N. Malpica, C. de Solórzano, J. Vaquero, A. Santos, I. Vallcorba, J. Garcia-Sagredo, and F. del Pozo, "Applying watershed algorithms to the segmentation of clustered nuclei," *Cytometry Part A*, vol. 28, no. 4, pp. 289-297, 1997.
- [89] L. True, and X. Gao, "Quantum dots for molecular pathology: their time has arrived," *Journal of Molecular Diagnostics*, vol. 9, no. 1, pp. 7, 2007.
- [90] A. Young, M. Amin, J. Petros, M. Natan, S. Nie, and M. Wang, "Nanomolecular histopathology for renal tumor classification." pp. 723-726.
- [91] S. Kothari, Q. Chaudry, and M. Wang, "Automated cell counting and cluster segmentation using concavity detection and ellipse fitting techniques." pp. 795-798.
- [92] M. Hall-Beyer, "The GLCM tutorial home page," p. 26.
- [93] J. R. Hand, and A. C. Broders, "Carcinoma of the kidney; the degree of malignancy in relation to factors bearing on prognosis," *The Journal of urology: official journal of the American Urological Association, Inc*, pp. 199, 1932.
- [94] N. S. Goldstein, "The current state of renal cell carcinoma grading," *Cancer*, vol. 80, no. 5, pp. 977-980, 1997.
- [95] P. Hermanek, A. Sigel, and S. Chlephas, "Histological grading of renal cell carcinoma," *European urology*, vol. 2, no. 4, pp. 189, 1976.
- [96] B. Delahunt, and J. N. Nacey, "Renal cell carcinoma II. Histological indicators of prognosis," *Pathology*, vol. 19, no. 3, pp. 258-263, 1987.
- [97] R. Montironi, A. Santinelli, R. Pomante, R. Mazzucchelli, P. Colanzi, A. Longatto Filho, and M. Scarpelli, "Morphometric index of adult renal cell carcinoma," *Virchows Archiv*, vol. 437, no. 1, pp. 82-89, 2000.
- [98] B. Delahunt, P. B. Bethwaite, J. N. Nacey, and J. L. Ribas, "Proliferating cell nuclear antigen (PCNA) expression as a prognostic indicator for renal cell carcinoma: Comparison with tumour grade, mitotic index, and silver-staining nucleolar organizer region numbers," *The Journal of pathology*, vol. 170, no. 4, pp. 471-477, 2005.
- [99] S. A. Fuhrman, L. C. Lasky, and C. Limas, "Prognostic significance of morphologic parameters in renal cell carcinoma," *The American Journal of Surgical Pathology*, vol. 6, no. 7, pp. 655, 1982.
- [100] S. J. K. Pedersen. "Circular Hough Transform," nov 18 2010, 2010; [http://www.cvmt.dk/education/teaching/e07/MED3/IP/Simon\\_Pedersen\\_CircularHoughTransform.pdf](http://www.cvmt.dk/education/teaching/e07/MED3/IP/Simon_Pedersen_CircularHoughTransform.pdf)
- [101] F. Chiarugi, S. Colantonio, D. Emmanouilidou, D. Moroni, F. Perticone, A. Sciacqua, and O. Salvetti, "ECG and echocardiography processing for decision support in heart failure." pp. 649-652.

## VITA

### **QAISER M. CHAUDRY**



QAISER CHAUDRY was born in Rawalpindi, Pakistan. He received his B.E. from College of Aeronautical Engineering, Risalpur, Pakistan and M.S. from Center for Advance Studies in Engineering, Islamabad, Pakistan in 2004. He served in Pakistan Army and worked with Microsoft Research. He carried out this research as a graduate student at School of Electrical and Computer Engineering, Georgia Institute of Technology. His primary research area is image analysis and bio-informatics with special focus to analysis and quantification of biomarker signatures in pathological images.