

University of Louisville

## ThinkIR: The University of Louisville's Institutional Repository

---

Electronic Theses and Dissertations

---

5-2016

### Automated analysis of cancerous histological samples.

Thomas Allen Tennill  
*University of Louisville*

Follow this and additional works at: <https://ir.library.louisville.edu/etd>



Part of the [Biomedical Engineering and Bioengineering Commons](#)

---

#### Recommended Citation

Tennill, Thomas Allen, "Automated analysis of cancerous histological samples." (2016). *Electronic Theses and Dissertations*. Paper 2371.

<https://doi.org/10.18297/etd/2371>

This Master's Thesis is brought to you for free and open access by ThinkIR: The University of Louisville's Institutional Repository. It has been accepted for inclusion in Electronic Theses and Dissertations by an authorized administrator of ThinkIR: The University of Louisville's Institutional Repository. This title appears here courtesy of the author, who has retained all other copyrights. For more information, please contact [thinkir@louisville.edu](mailto:thinkir@louisville.edu).

AUTOMATED ANALYSIS OF CANCEROUS HISTOLOGICAL SAMPLES

By

Thomas Allen Tennill  
B.S., University of Louisville, 2014

A Thesis  
Submitted to the Faculty of the  
University of Louisville  
J.B. Speed School of Engineering  
As Partial Fulfillment of the Requirements  
For the Professional Degree

MASTER OF ENGINEERING

Department of Bioengineering

May 2016



AUTOMATED ANALYSIS OF CANCEROUS HISTOLOGICAL SAMPLES

Submitted by: \_\_\_\_\_  
Thomas Allen Tennill

A Thesis Approved On

\_\_\_\_\_  
(Date)

by the Following Reading and Examination Committee:

\_\_\_\_\_  
Hermann Frieboes, Thesis Director

\_\_\_\_\_  
Martin O'Toole

\_\_\_\_\_  
Stuart Williams

## ACKNOWLEDGEMENTS

A special thank you to Dr. Hermann Frieboes for his continued guidance and mentorship throughout my academic career. An additional thank you to Dr. Mitchell Gross for his assistance and input on this work. Lastly, thank you to the members of my committee: Dr. Martin O'Toole and Dr. Stuart Williams.

## ABSTRACT

The use of immunohistochemistry has become commonplace in the field of cancer diagnosis. A major limitation to the ability to use this tool effectively is the time consuming tasks required for the analysis of data due to the fact that it is largely done by hand. Additionally, because of this, there is an inherent level of subjectivity in the results obtained from this process that may depend on who it is conducting the analysis. Therefore, there exists a need for a method that is able to quantify results from immunohistochemical techniques in a way that it is both time-effective and consistent in how each sample is treated.

In this study a program was developed that was able to give a quantitative analysis of DAB stained prostate cancer samples that mimics the results obtained by the conventional manual annotation method. This program was then used further to analyze much larger samples that would be too time consuming to analyze in the conventional way, as well as to analyze a large series of samples generated in a tissue microarray.

## TABLE OF CONTENTS

TITLE PAGE .....	i
APPROVAL PAGE .....	iii
ACKNOWLEDGEMENTS .....	iv
ABSTRACT .....	v
NOMENCLATURE .....	vii
LIST OF TABLES .....	viii
LIST OF FIGURES .....	ix
I.    INTRODUCTION .....	1
A. Background on Prostate Cancer .....	1
B. Background on Immunohistochemistry .....	3
C. Current Use of Automated Programs for Analysis of IHC Stained Samples .....	4
II.   INSTRUMENTATION AND EQUIPMENT .....	7
III.  METHODS .....	8
A. Tissue Micro Array .....	8
1. Program Startup .....	8
2. Reading in Image .....	9
3. Isolating Core Samples .....	9
4. Collecting/Managing Core Location Data .....	11
5. Reading in Large Scale Image/Translating Location Data .....	14
6. Processing of Large Scale Image .....	15
B. Whole Slice .....	19
1. Initial Setup .....	19
2. Image Segmentation .....	19
3. Section Processing .....	20
4. Data Management .....	21
C. Region of Interest .....	22
1. Initial Setup .....	22
2. Image Analysis .....	23
D. Processing Flow Chart .....	24
IV.  RESULTS .....	25
A. Comparing Program to Manual Results .....	25
B. ROI/Whole Slice/TMA Analysis .....	31
V.   DISCUSSION .....	39
VI.  CONCLUSIONS .....	41
VII. RECCOMENDATIONS .....	43
REFERENCES CITED .....	45

## NOMENCLATURE

PIN = Prostatic Intraepithelial Neoplasia

PIA = Proliferative Inflammatory Atrophy

IHC = Immunohistochemistry

TMA = Tissue Microarray

H&E = Hematoxylin and Eosin

DAB = 3,3'-Diaminobenzidine

GHz = Gigahertz

CPU = Central Processing Unit

ROI = Region of Interest

StDev = Standard Deviation

t = T-Score

$\bar{x}_1, \bar{x}_2$  = Sample Means

$S_1^2, S_2^2$  = Sample Variance

$N_1, N_2$  = Number of Samples

DoF = Degrees of Freedom

p = P-Value



## LIST OF TABLES

TABLE I – MiB-1 RESULTS .....	33
TABLE II – CD31 RESULTS .....	34
TABLE III – MiB-1/CD31 RATIO RESULTS .....	34
TABLE IV – ROI STATISTICAL SIGNIFICANCE .....	36
TABLE V – WHOLE SLICE STATISTICAL SIGNIFICANCE .....	37
TABLE VI – TMA STATISTICAL SIGNIFICANCE .....	38

LIST OF FIGURES

FIGURE 1 – Core Isolation .....10

FIGURE 2 – Centroid Matrix Excerpt .....11

FIGURE 3 – Location Matrix/TMA Direct Comparison .....13

FIGURE 4 – Core Numbering Adjustment .....14

FIGURE 5 – Translation of Location Data to Large Scale Image.....15

FIGURE 6 – Processing of TMA Core.....17

FIGURE 7 – TMA Data Exportation.....18

FIGURE 8 – Whole Slice Image Segmentation .....20

FIGURE 9 – Processing Flow Chart.....24

FIGURE 10 – ROI MiB-1 Manual/Program Comparison .....26

FIGURE 11 – ROI CD31 Manual/Program Comparison .....27

FIGURE 12 – ROI MiB-1/CD31 Ratio Manual/Program Comparison.....28

FIGURE 13 – Weighted Program and Manual MiB-1 Results .....29

FIGURE 14 – Weighted Program and Manual CD31 Results .....29

FIGURE 15 – Weighted MiB-1 Results with Error Bars .....30

FIGURE 16 – Weighted CD31 Results with Error Bars .....30

FIGURE 17 – MiB-1 Results.....32

FIGURE 18 – CD31 Results.....33

FIGURE 19 – MiB-1/CD31 Ratio Results .....34

## I. INTRODUCTION

### A. Background on Prostate Cancer

The prostate is a gland present only in males that sits between the urinary bladder and the rectum whose job is to make a fluid that nourishes sperm cells which make the semen more liquid (1). Prostate cancer is the second most common form of cancer in American men, affecting 1 out of 7 men in their lifetime, and the second leading cause of cancer death in American men, killing 1 in 38 American men. It is estimated to create 220,800 new cases and 27,540 deaths in just a year. The risk of prostate cancer increases with age with 60% of cases being in men age 65 and older and the average age of the patient at diagnosis being 66(2). Almost all cases of prostate cancer are adenocarcinomas, a type of cancer that starts in the gland cells which, in the prostate, create the prostate fluid. There are other, rarer forms of prostate cancer which include sarcomas, small cell carcinomas, neuroendocrine tumors, and transitional cell

carcinomas. Prostate cancer might start as a pre-cancerous condition, one example of which is Prostatic Intraepithelial Neoplasia (PIN). PIN involves prostate gland cells that appear abnormal under a microscope but do not grow into other parts of the prostate. PIN is extremely common as almost 50% of men have low-grade PIN by age 50. Individuals possessing high-grade PIN have about a 1 in 5 chance of also having prostate cancer present in another part of the prostate. Another pre-cancerous condition that may lead to prostate cancer Proliferative Inflammatory Atrophy (PIA). PIA involves prostate cells that look smaller than normal under microscope after a biopsy while also possessing signs of inflammation. PIA is believed to lead to either high-grade PIN or to prostate cancer directly (1).

There are various risk factors for prostate cancer including age, race/ethnicity, geography, family history, genetic mutations, diet, and undergoing a vasectomy. Prostate cancer is rare in men under age 40, however the risk rises rapidly after age 50. Prostate cancer is most common in African-American men and Caribbean men of African ancestry, with these individuals being twice as likely to die from prostate cancer as white men. Asian-American and Hispanic/Latino men are less likely to be diagnosed with prostate cancer than non-Hispanic white men. Prostate cancer is most prevalent in North America, northwestern Europe, Australia, and the Caribbean Islands. Having either a father or brother who has had prostate cancer more than doubles the risk of the individual developing prostate cancer. Increased risk for prostate cancer has also been found to be tied to specific genetic mutations on the BRCA1 and BRCA2 genes as well as individuals with Lynch Syndrome. Eating a diet consisting of red meat, high-fat dairy products, along with a lack of fruits and vegetables can slightly increase the risk of developing

prostate cancer. Also, men who have had a vasectomy performed on them have a slightly increased risk for developing prostate cancer (3). There are a variety of treatment options available to patients with prostate cancer including active surveillance, surgery, radiation therapy, cryotherapy, hormone therapy, chemotherapy, vaccine treatments, and bone-directed treatments. The factors affecting which treatment is best for the patient include the age and life expectancy of the patient, other health conditions, stage/grade of the cancer, necessity for immediate treatment, success rates for different treatments, and risk of possible side effects of these treatments (4).

### B. Background on Immunohistochemistry

Immunohistochemistry has been used as a complimentary method in the diagnosis of various forms of cancer. Specific tumor markers, primarily proteins, have been identified whose levels can be used as signals for the presence of tumor cells (6). CD31 (aka platelet-endothelial cell adhesion molecule type 1) is an immunohistochemical stain that indicates endothelial cells, granulocytes, monocytes, and platelets. CD31 is primarily used to identify tumors that are of endothelial origin. MiB-1 is a nuclear non-histone protein that is present in all stages of the cell cycle except G0. Constantly proliferating cells express this protein and, therefore, it is useful in estimating the growth fraction of both benign and malignant tissue (5).

The analysis of immunohistochemically stained tissue samples has traditionally been done manually by a pathologist through a process that is both time consuming and subjective. The use of tissue micro arrays (TMA's) for many studies can result in the development of hundreds of samples needing to be analyzed at a time. This has led to a need for an automated process that can reduce the time necessary for sample analysis and perform consistently for every sample.

### C. Current Use of Automated Programs for Analysis of IHC Stained Samples

Automated programs have been developed to meet the need for quick analysis tools in analyzing IHC stained tissue in specific applications. One example of the use of automated programs is to separate areas of ovarian carcinoma from the remainder of the sample (8). This program was able to successfully segregate the portions of various TMA cores that contained carcinomas from the remainder of the sample and the background using the existing Genie Histology Pattern Recognition System (Leica Biosystems). This system, however was not able to be used in a fully automated way as it still required a pathologist and technician to identify tumor regions in samples in order to create the input parameters for the software and also requires multiple repetitions of the training algorithm before it was able to produce results mimicking the original work done by the pathologist and technician, with each iteration of the training program requiring the user to provide feedback and make adjustments to the program's results.

Another main limitation driving the need for the long calibration periods of the current methods available is that each IHC stain and cancer type requires individual attention for method development. While this system is robust in its ability to potentially be used for multiple different types of stain and tissue, the calibration methods can lead to a less than ideal experience for those needing to routinely run high throughput analysis on samples of consistent tissue and staining methods.

Another example of using computers to analyze IHC staining is CRImage, which has been used in the analysis of hematoxylin and eosin (H&E) staining (9, 10). CRImage is a package extension of the R programming language, which can be downloaded and used freely. This program is capable of classifying cells, segmenting samples, and calculating tumor cellularity, on top of R's native statistical analysis capabilities. There are, however limitations involved with CRImage that can be improved upon. Firstly, in order to use CRImage, you must have some level of knowledge of the R programming language due to the lack of a graphical user interface which can present a large obstacle for medical doctors who would like to make use of the program. Additionally, CRImage is not a fully automated analysis system. While it is capable of image segmentation and sample classification without the need for user input, it still requires the user to enter commands for every process they want to use on the image and each image must be processed individually. CRImage is also incapable of being used on Aperio Image Slides, a popular format for saving high resolution pathological images, due to their use of the SVS file format. Workarounds for this problem exist but require the user to change the format of the image outside of the CRImage package. Lastly, the use of the CRImage analysis package is limited to samples possessing H&E stain so this is not a



practical solution on its own for those interested in processing slides from various stain types.

Computer aided techniques have also been used extensively for cell counting (11, 12, 13). MATLAB has been a successful platform for developing multiple cell counting algorithms. These programs tend to have a more limited applications than the previous options, providing only information on how many cells are present and not analyzing the characteristics of those cells. Many of these programs are intended for use with either dark field microscopy or grayscale images that possess high contrast between the items of interest and the background, making separating cells from the background of the image simpler. These programs would not inherently be usable with DAB staining methods but indicate that MATLAB is a viable platform for biological image analysis software.

## II. INSTRUMENTATION AND EQUIPMENT

This program was designed using Matlab version 8.3 (R2014a) on a laptop running Windows 8.1 (64 bit). The laptop possessed 8 gigabytes of physical memory and an Intel 1.80 GHz Core i7-4500u CPU.

### III. METHODS

#### A. Tissue Micro Array

##### 1. Program Startup

Due to the large size of the files for the high resolution images, it is possible that the computer may run out of memory. For this reason, the program begins by closing all open windows within MATLAB, clearing the command screen, and clearing all variables stored in memory in order to avoid limitations associated with the hardware of the computer it is being run on. Additionally, warnings have been turned off to prevent warnings associated with displaying the images.

## 2. Reading in Image

A scaled down version of the TMA image is first loaded in to extract necessary information while conserving system memory. Information regarding the size of this image is gathered in order to create an all-black image of the same size.

## 3. Isolating Core Samples

The first step of separating the cores on the TMA's is to remove the background of the image. Each pixel is examined to determine whether it is part of a core or the background. If the pixel is determined to be part of the background then the corresponding pixel on the black image is changed to white. Some samples are not continuous, therefore it is necessary to manipulate the image to merge those sections. This is accomplished by creating a structuring element, using it to erode/dilate the image, creating a perimeter around the outside of the object, and filling all of the holes within the object to leave one whole area representing the location of the core. In order to remove unwanted objects and samples that have merged together from the TMA the areas of the objects are analyzed, turning all areas that are either exceedingly small or large back to black. A visual representation of this process can be seen below:

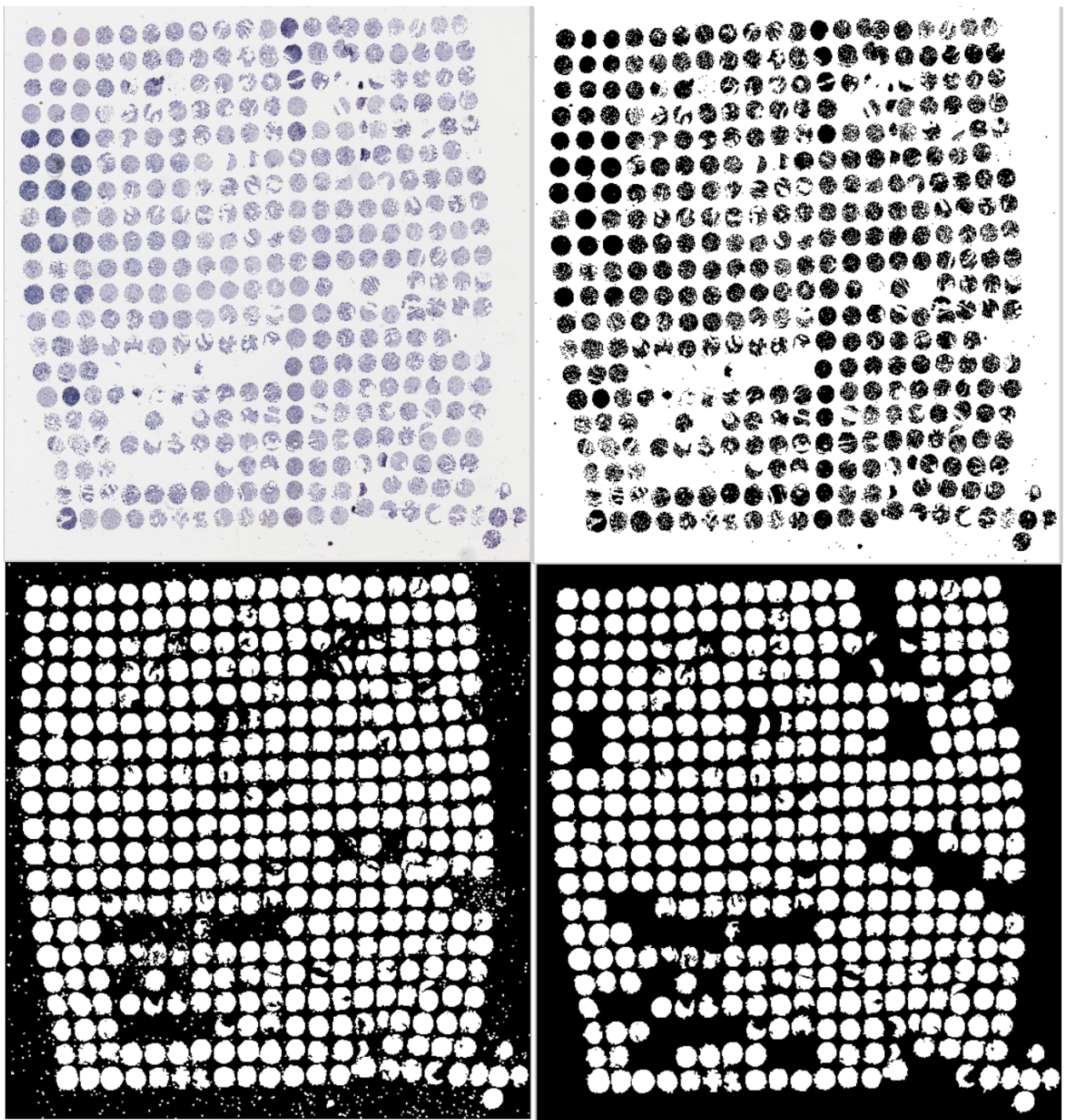


FIGURE 1 – Core Isolation

Top Left: The original image is shown. Top Right: Background is separated from the samples. Bottom Left: Holes in the samples are filled to create solid objects. Bottom Right: Small objects are removed as well as samples that are overlapping each other.

#### 4. Collecting/Managing Core Location Data

Centroid locations of all of the remaining white objects are collected and exported to a matrix with each row representing a different core, the first column representing the horizontal position of the centroid, and the second column representing the vertical position of the centroid. A third column is then created and filled with the value equal to the sum of the first two columns and the entire matrix is sorted by the third column to find the core closest to the image's origin, i.e. the top left corner of the image.

255	282	537
446	286	732
240	494	734
632	284	916
432	491	923
233	694	927
816	271	1087
629	487	1116
436	695	1131
220	912	1132

FIGURE 2 – Centroid Matrix Excerpt

Excerpt from the centroid matrix. The left column shows the X position of the centroid, central column shows the Y position of the centroid, and the right column shows the sum of the first two columns.

A separate location matrix is created with size equal to that of the grid of cores present on the TMA. The first entry in the centroid matrix is then labeled as the first entry, (1,1), in the location matrix. The program then references the centroid location of

the original core to search for adjacent cores, focusing first on completing the first column by searching the centroid matrix for entries with similar horizontal position to the previous core and vertical positions within a specified range away from the previous core. If a centroid is found to be in this expected location then that core is referenced in the location matrix by filling its relative position with the row number of its entry in the centroid matrix and the process is repeated to find the next core along the column. If no core is found for an expected area then the search is expanded to look for a centroid that would be located two places on the grid away from the previous core and the search is expanded in this fashion until the next core is found. Once the first column is complete there is a check to make sure a proper amount of samples were found to create a baseline for the row information based off of that column. If the amount of samples is found to be insufficient then the program will instead attempt to create a baseline from the first row rather than the first column. Once the baseline row/column is established, the program will then fill in the columns/rows by similar means to the previous method, using the baseline row/column as the original reference.

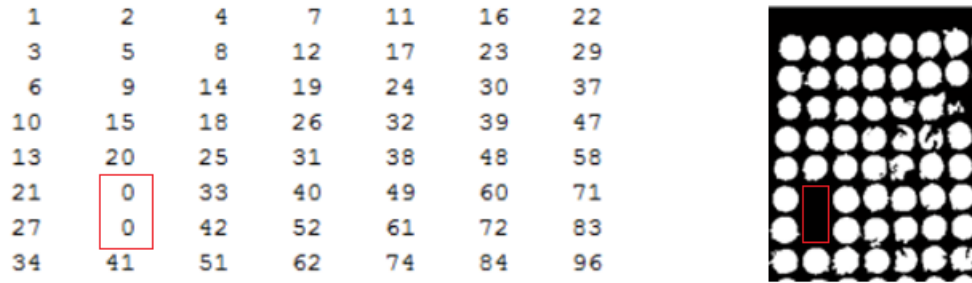


FIGURE 3 – Location Matrix/TMA Direct Comparison

Excerpt from the location matrix and the portion of the TMA it represents. The location matrix is filled in with the sample number in the corresponding location of the sample on the TMA. Locations where a core is not found are left as a 0, as highlighted by the red rectangles.

Once the location matrix is complete, two columns are added to the centroid matrix and are filled with the row and column information for where that sample was placed in the location matrix. This results in the assigning of the core to both a column and row for future reference. Due to the fact that the TMA images are oriented opposite of the conventional numbering system being used for the rows and columns of the samples the value of the row and column assigned to each sample is inverted so that samples assigned to the first row or column are now labeled as being in the last row or column and vice versa.



255	282	537	1	1	255	282	537	20	20
446	286	732	2	1	446	286	732	19	20
240	494	734	1	2	240	494	734	20	19
632	284	916	3	1	632	284	916	18	20
432	491	923	2	2	432	491	923	19	19
233	694	927	1	3	233	694	927	20	18
816	271	1087	4	1	816	271	1087	17	20
629	487	1116	3	2	629	487	1116	18	19
436	695	1131	2	3	436	695	1131	19	18
220	912	1132	1	4	220	912	1132	20	17

FIGURE 4 – Core Numbering Adjustment

Left of red line: Centroid matrix updated with rankings based on core locations. The first number represents which column the core is in while the second number represents which row the core is in. Right of red line: Centroid matrix once the row and column information has been adjusted to match the conventional numbering system.

#### 5. Reading in Large Scale Image/Translating Location Data

The large scale TMA image is now read into the program so that the cores can be analyzed at a higher resolution. Due to the fact that the location information gathered for the centroid locations does not correspond to their locations on the larger image the pixel locations must be multiplied by a scale factor equal to the scale factor of the image sizes.

255	282	537	20	20	1020	1128	537	20	20
446	286	732	19	20	1784	1144	732	19	20
240	494	734	20	19	960	1976	734	20	19
632	284	916	18	20	2528	1136	916	18	20
432	491	923	19	19	1728	1964	923	19	19
233	694	927	20	18	932	2776	927	20	18
816	271	1087	17	20	3264	1084	1087	17	20
629	487	1116	18	19	2516	1948	1116	18	19
436	695	1131	19	18	1744	2780	1131	19	18
220	912	1132	20	17	880	3648	1132	20	17

FIGURE 5 – Translation of Location Data to Large Scale Image

Left: excerpt from the original centroid matrix. Right: centroid matrix updated to accurately reflect centroid information on the higher resolution image. The first two columns, representing the X and Y positions of the center of the cores, have been scaled up to the new resolution.

Additionally, data is read in from a spreadsheet consisting of the row location, column location, and Gleason Scores of the cores on the TMA. This data is used to assign the proper Gleason Score to each core that was found by the program.

## 6. Processing of Large Scale Image

The scaled location data from the centroid matrix is used to identify the centroid locations on the large scale image. From the area surrounding the centroid location a composite image is generated that consists of only one sample. Another image of the same size as the single core image is created to document the analysis of the core. First, the core is analyzed pixel by pixel to determine whether each pixel is part of the

background or the sample in question. If it is determined that the pixel is part of the background then the corresponding pixel in the analysis image is set to black, whereas if the pixel is determined to be part of the sample then the corresponding pixel is set to white. The program now takes the sum total of all of the black pixels in the analysis image to document how much of the image is background and the sum total of all of the white pixels to determine how much of the image is part of the sample. Next the core is analyzed pixel by pixel to determine if there is brown stain present in that pixel. If it is determined that there is brown stain present in the pixel then the corresponding pixel in the analysis image is set to black, whereas if there is no brown stain present in the pixel the corresponding location is set to white. The total amount of stain present in the sample is found by taking the sum total of all of the black pixels in the analysis image.

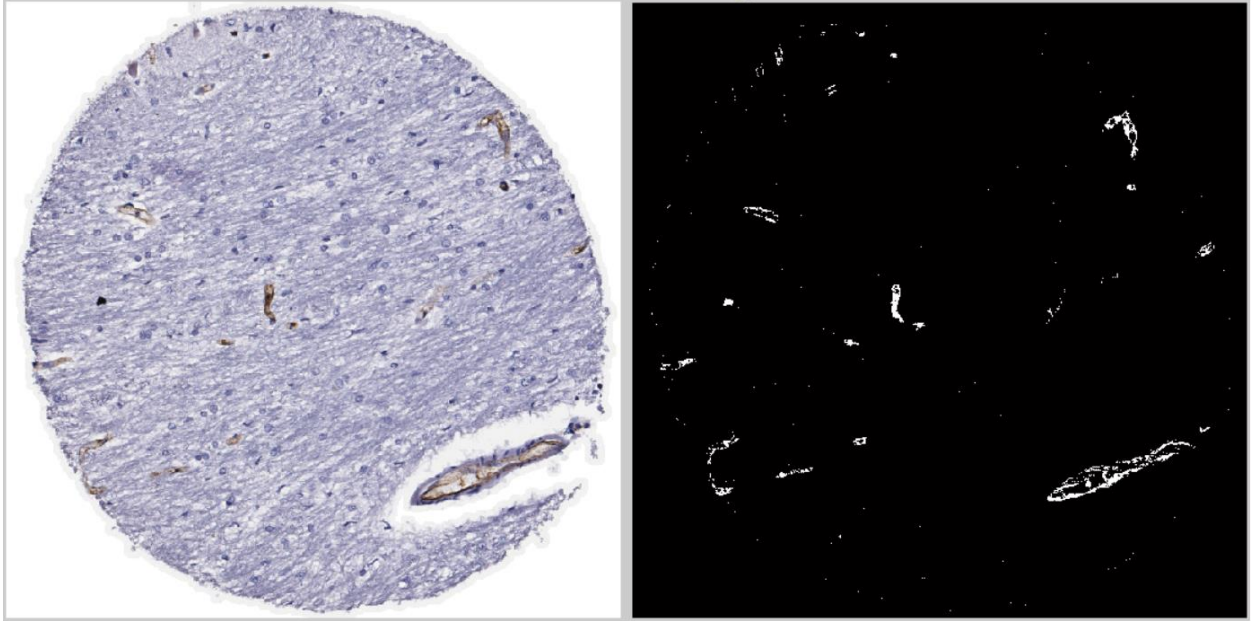


FIGURE 6 – Processing of TMA Core

Left: Sample of a core being analyzed individually. Right: core on the left after being analyzed for the presence of stain. The white areas reflect the location of the brown stain in the left image.

From here, the percent of the core that is stained is found by dividing the number of pixels found to possess the brown stain by the number of pixels that were not found to be part of the background and multiplying by 100. This data is then recorded in a new column of the centroid matrix and the process is repeated on the next core until all have been analyzed. Once all core analysis is complete, the centroid matrix, consisting of core centroid locations, row and column of the core, Gleason Score of the core, and the percentage of the core that possessed the brown stain, is exported and saved in an excel file for further analysis.

	A	B	C	D	E	F	G
1	1020	1128	537	20	20	0.053254	
2	1784	1144	732	19	20	0.124575	
3	960	1976	734	20	19	0.039264	
4	2528	1136	916	18	20	0.044647	
5	1728	1964	923	19	19	0.076324	
6	932	2776	927	20	18	0.115486	
7	3264	1084	1087	17	20	0.04277	6
8	2516	1948	1116	18	19	0.069113	
9	1744	2780	1131	19	18	0.131296	
10	880	3648	1132	20	17	0.592754	

FIGURE 7 – TMA Data Exportation

Excerpt from exported excel file. The first two columns represent the X and Y locations for the center of the core. The third column represents distance from the image's origin. Columns 4 and 5 represent which column and row the core is located in. The sixth column is the percentage of the core that was found to possess the stain, while the seventh column indicates the Gleason Score of the core. Several samples are used as controls on the TMA, therefore not all samples have a Gleason Score associated with them.

## B. Whole Slice

### 1. Initial Setup

The processing of whole slice images allows for the user to process a series of images automatically. In order for the program to do so, the user must set up a folder including all images that are to be processed. Upon program startup the user will be presented with a file explorer which they will use to select the folder in which they've put the images. The program then reads the number of SVS images present in the folder to determine how many images will need to be analyzed.

### 2. Image Segmentation

The program begins by accessing the file name of the image currently being analyzed for labeling purposes. The program then reads in that image and collects both the height and width of the image. Due to limitations associated with both the size of the images and the hardware used to process them, each image of the whole tissue slice needs to be segmented to be processed one piece at a time. The program divides the original image up into a four by four grid, creating sixteen pieces equal in size, as seen below.

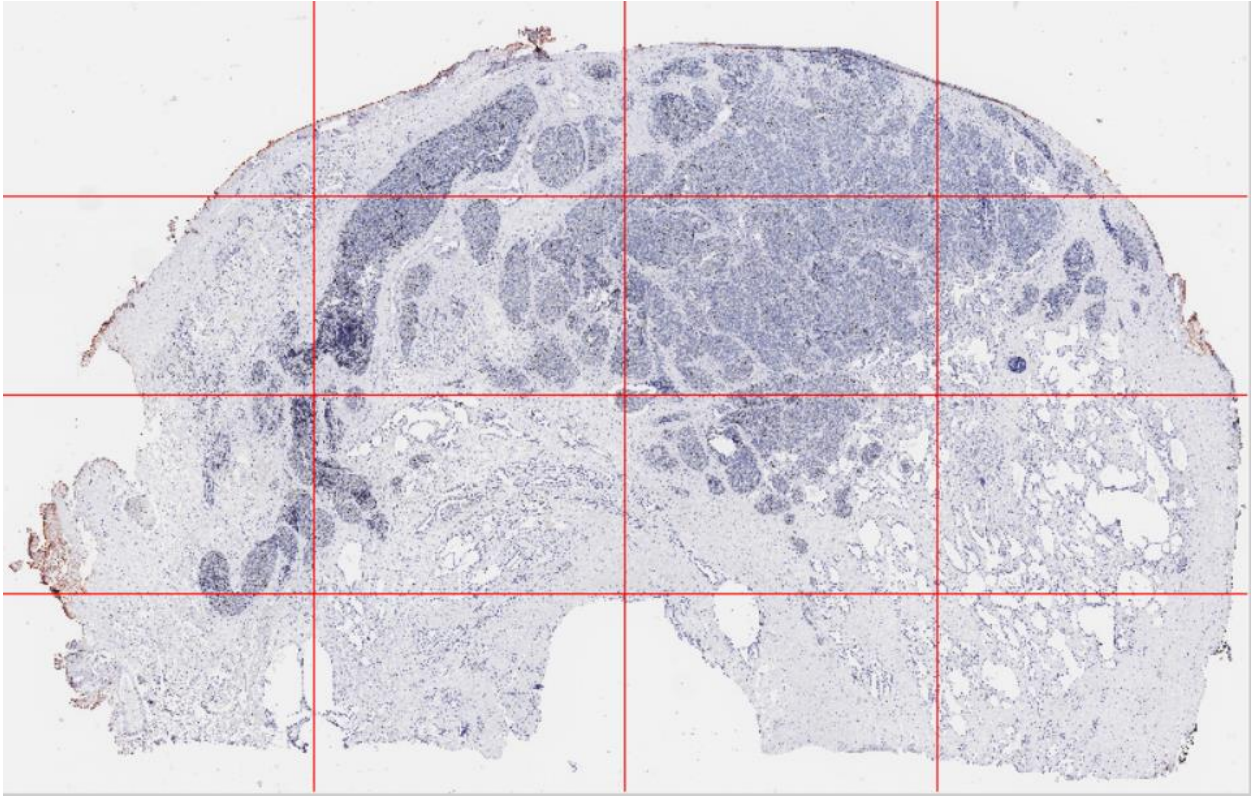


FIGURE 8 – Whole Slice Image Segmentation

The program then creates a separate image of size equal to that of the sixteen sections to be used during the analysis. The program also creates variables that will be used to store values for the amount of pixels in the image that represent the background, sample, stained sample, and non-stained sample and sets them all equal to zero.

### 3. Section Processing

A series of nested for-loops begins to have the program begin to process all sixteen pieces one at a time. To reduce strain on the hardware running the program, each

sections analysis overwrites the previous. The program first analyzes each pixel within the section by categorizing it to either be a part of the sample or the background. The corresponding pixel location in the composite image is then changed to either be black or white based on the category. The total number of black pixels in the composite image is calculated, as well as the total number of white pixels in the image, and then the number of pixels found to be in each category is added to the variable associated with the category outside the for-loop during image segmentation. Next each pixel of the section is categorized to either be containing the brown stain or not in a method similar to the sample/background categorization. The number of white pixels present in the composite image is then counted to be representative of the amount of stained sample in this section. The amount of non-stained sample in this section is calculated by subtracting the amount of stained sample from the amount of total sample for the section. These values are then added to the variable associated with the category outside the for-loop for storage. This process then continues for each of the remaining 16 sections, continually adding the values found for each section to the cumulative total being tracked.

#### 4. Data Management

Once the totals for the entire image has been calculated, the total pixels present in the image is calculated by adding the amount of background, stained sample, and non-stained sample pixels together. This value is then used to find the percent of the total image that is represented by background, stained sample, and non-stained sample. The



proportion of the tissue that possesses the brown stain is found by dividing the number of pixels found to possess the stain by the total amount of pixels making up the sample. At the end of the processing of each image the program exports the file name, percentage of the image that is background, possesses brown stain, and normal tissue, and the proportion of the tissue that possesses the brown stain to a new row in an Excel file for further analysis.

### C. Regions of Interest

#### 1. Initial Setup

The processing of the regions of interest allows for the continuous processing of multiple images provided that they are all located within the same folder. The program begins by prompting the user to select the folder containing the images of interest. The program then reads the number of tiff images present in the folder to determine how many times the analysis needs to be run.

## 2. Image Analysis

For each image, the program begins by reading in its filename to label the sample. Next, the program collects size information about the image to create a blank composite image of duplicate height and width. Due to the relatively small file sizes associated with the ROI's, as opposed to the whole slice images and the TMA's, the entire image can be processed in one piece. The first step is to analyze the image to separate background from the sample. As the program goes through the ROI pixel by pixel, it sets the corresponding pixel in the composite image to either be black if it is a background pixel, or white if it is a part of the sample. The sum total of all of the white pixels is then taken to determine the area occupied by the sample. The program then goes over the original image again to separate pixels containing the stain from the rest of the image. If the program determines the pixel contains the stain it will turn the corresponding location in the composite image white, while leaving the location black if it does not possess the stain. The total number of white pixels in the composite image is then calculated to find the area occupied by the stain and this value is used to calculate what percent of the sample present in the image contains the stain. The program then exports the images filename and this stain percentage to an excel file before beginning this process again on the next image.

D. Processing Flow Chart

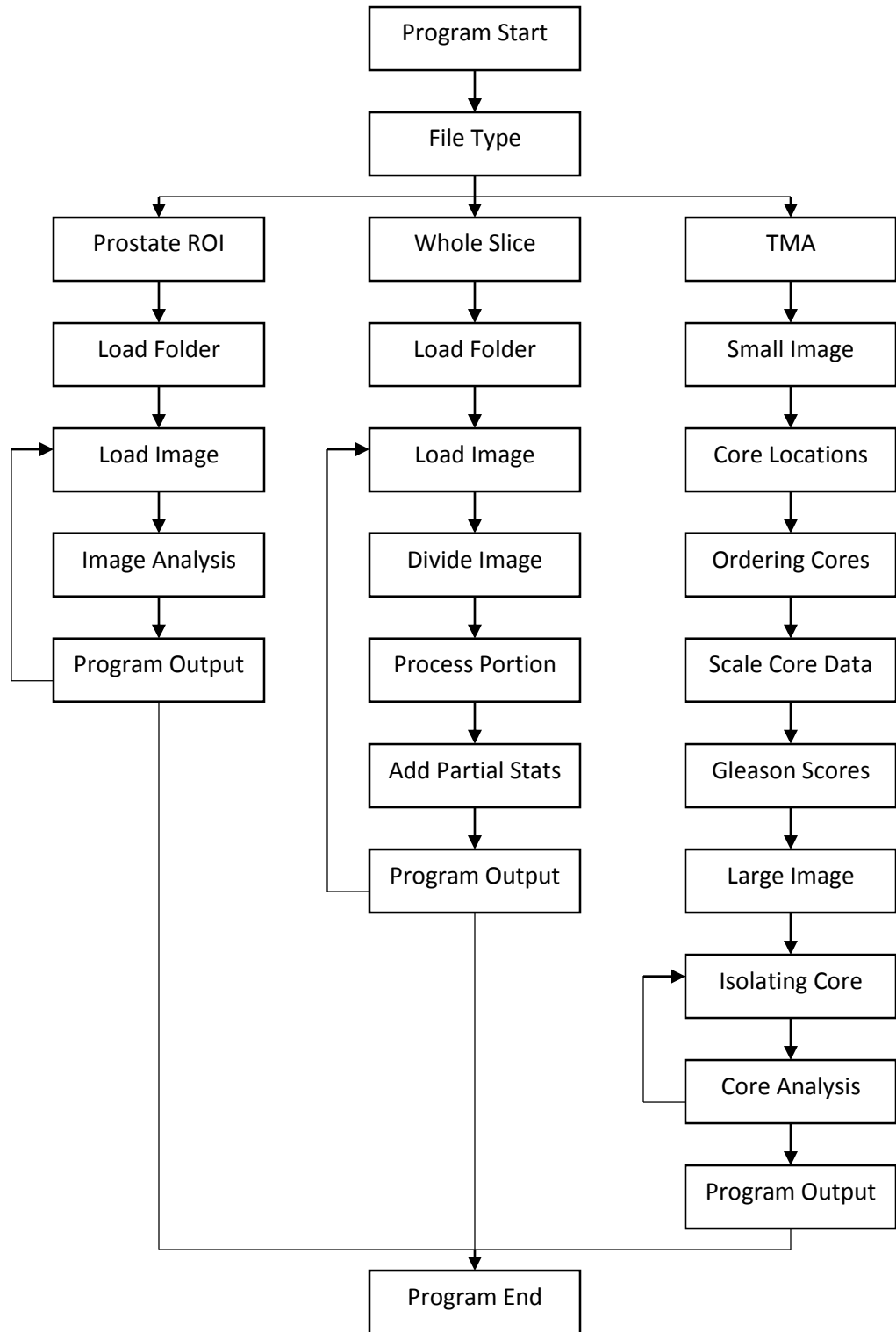


FIGURE 9 – Processing Flow Chart

## IV. RESULTS

### A. Comparing Program to Manual Results

The process of designing an automated program began with designing a system that could replicate the current process. The current process had been used to analyze various regions of interest from existing tissue slices so these were used as the benchmark. The program ran analysis for these regions possessing both MiB-1 and CD31 stains and calculated the ratio of the two for each region. For both the MiB-1 and CD31 stains the program underreported in comparison to the manual method but followed the same general trend between samples. This resulted in the ratio of the MiB-1 stain to the CD31 stain being consistent with the data collected manually.

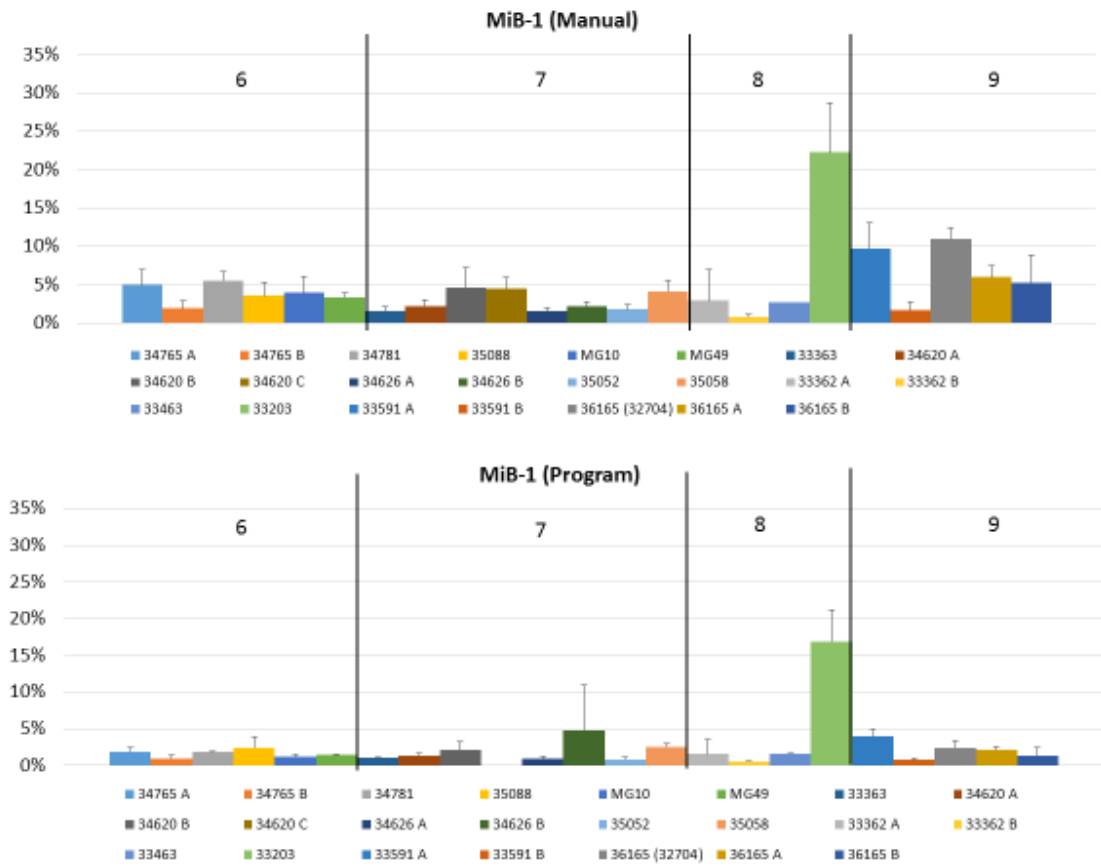


FIGURE 10 – ROI MiB-1 Manual/Program Comparison

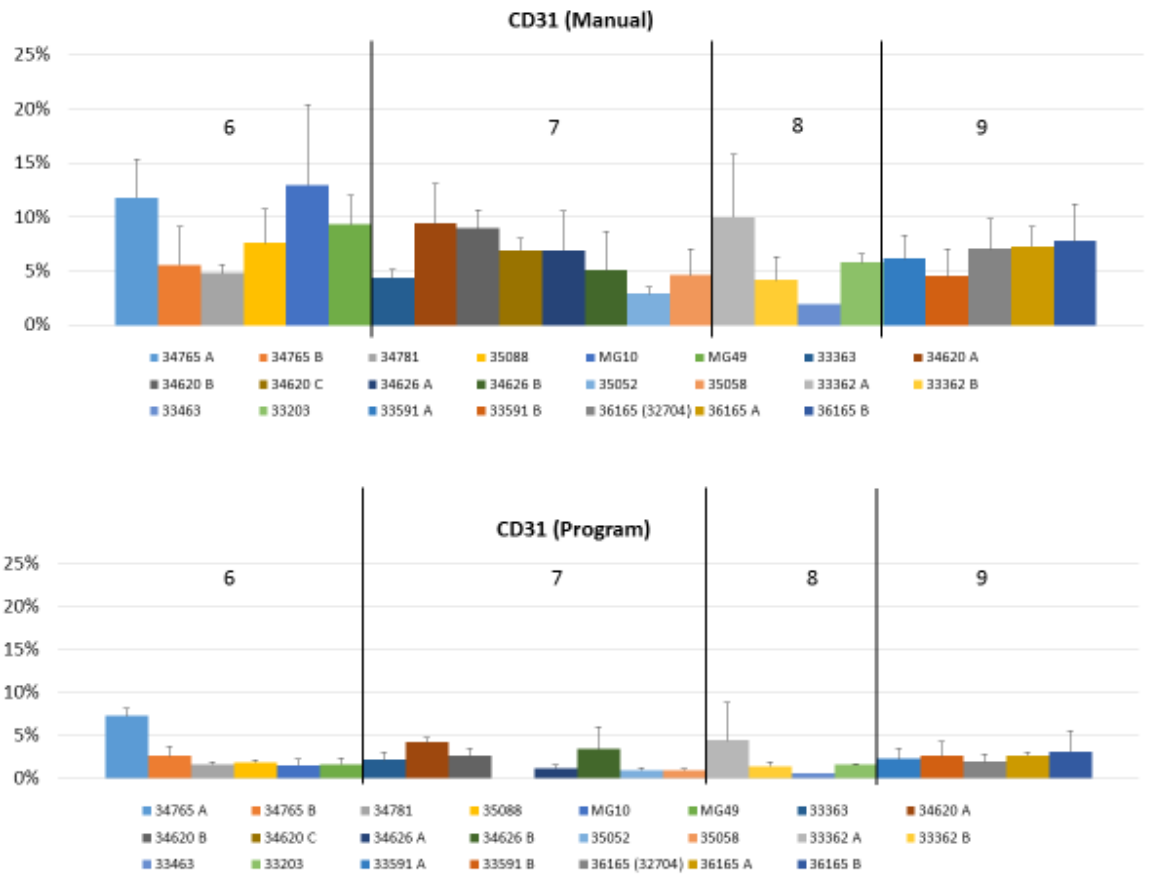


FIGURE 11 – ROI CD31 Manual/Program Comparison

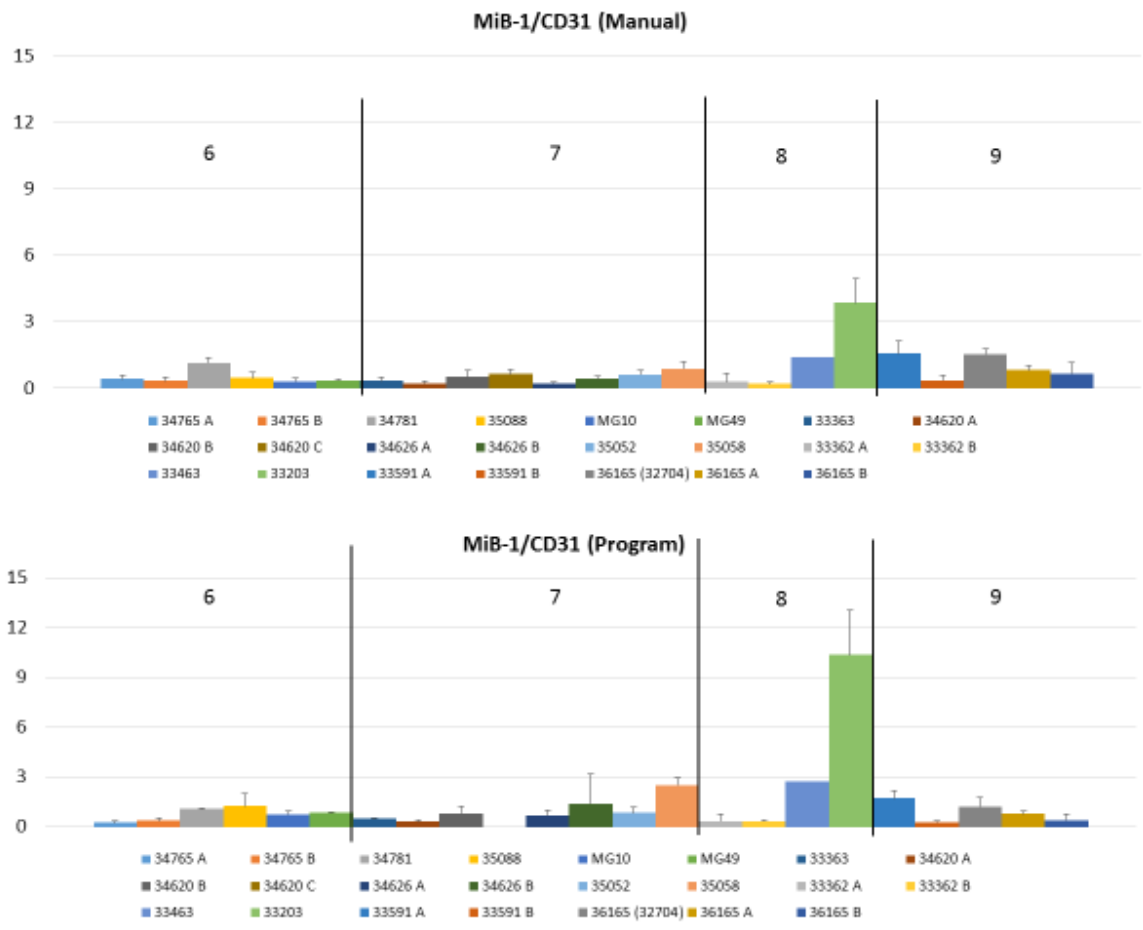


FIGURE12 – ROI MiB-1/CD31 Ratio Manual/Program Comparison

In order to show that the program was showing similar results as the manual process the data was normalized in respect to its method. To do this, the reported average stain concentrations for each sample were divided by the sample reported to possess the highest stain concentration, creating a scale from zero to one. The standard deviations measured from the method were also transformed to match this scale. The resultant data was then plotted to visualize the comparison.

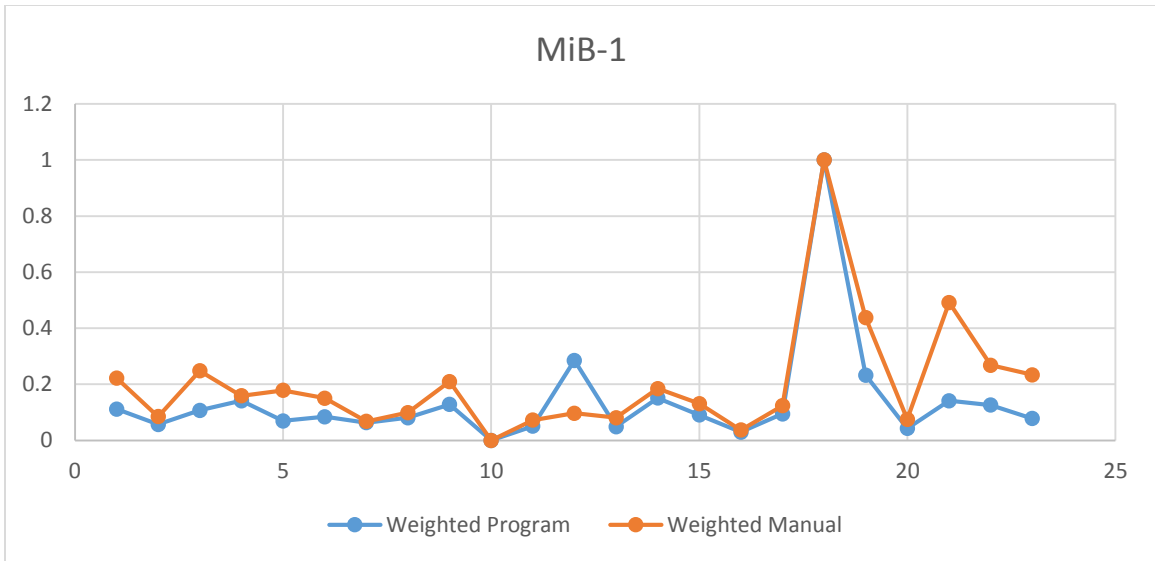


FIGURE 13 – Weighted Program and Manual MiB-1 Results

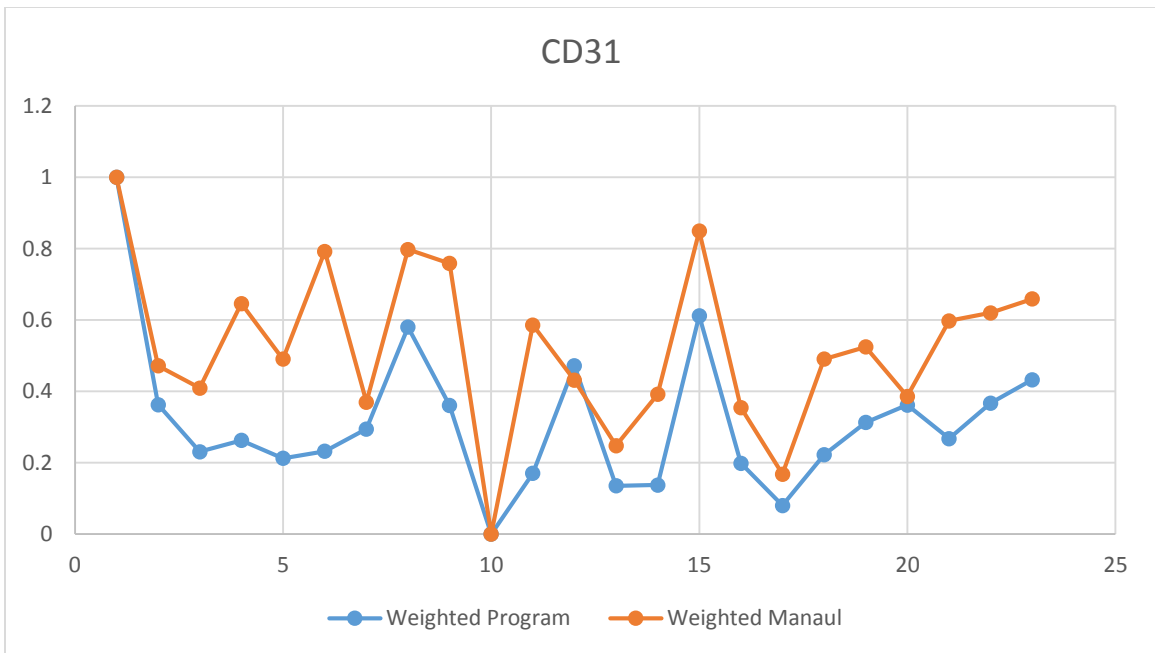


FIGURE 14 – Weighted Program and Manual CD31 Results

Visual inspection seemed to indicate that the program was reporting the same relative results as the manual method due to the similarity in the trends of the graph. To confirm this the standard deviations were added to the plots as error bars. In doing this,



having an overlap of the error bars for the two points representing the same sample would indicate that the weighted values are not distinctly different.

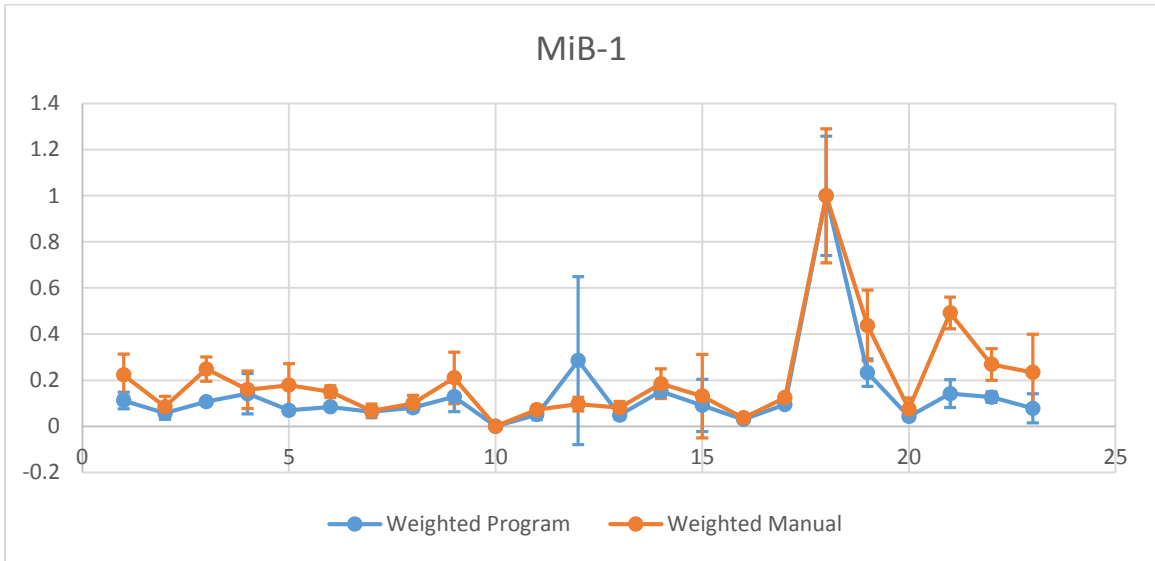


FIGURE 15 – Weighted MiB-1 Results with Error Bars

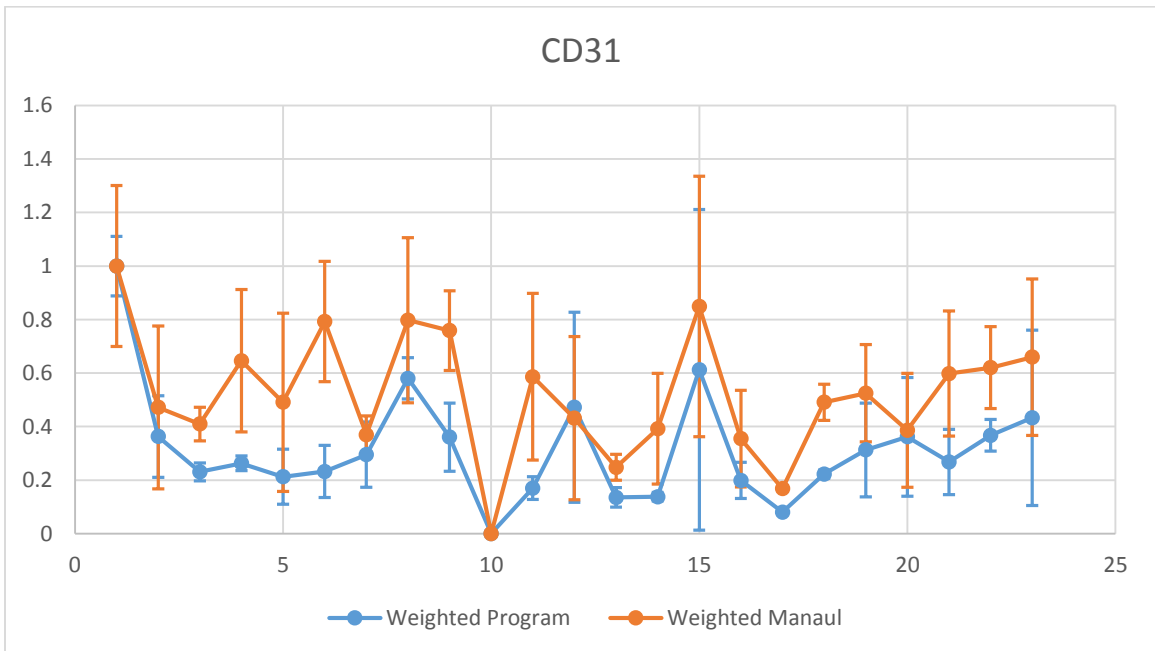


FIGURE 16 – Weighted CD31 Results with Error Bars

Analysis of this technique shows that 30 of the 43 data points (~70%) include overlap between the ranges of (reported averages  $\pm$  one standard deviation) for the two methods. This provided sufficient support for the accuracy of the automated method to continue developing it for further analysis.

### B. ROI/Whole Slice/TMA Analysis

The program was used to analyze three different data sets: regions of interest, whole tissue slices, and tissue micro arrays. For each sample within each data set there was an image possessing MiB-1 stain and a separate image possessing CD31 stain. Once every sample had been analyzed the samples were separated into two categories based on the Gleason Score's for the patient, one group consisting of all samples that possessed Gleason Scores of either 6 or 7, the other consisting of all samples with Gleason Scores of 8 or 9. A t-test was then conducted on the difference between these two categories for each of the data sets in terms of their presence of MiB-1 stain, presence of CD31 stain, and the ratio of the amount MiB-1 stain the sample possessed to the amount of CD31 stain the sample possessed. The results of the t-test showed that the difference in the presence of the MiB-1 stain was only significant in the Region of Interest samples taken, the difference in the presence of CD31 stain was only significant in the TMA samples, while the difference in the MiB-1 to CD31 ratio was found to be significant in all data sets.

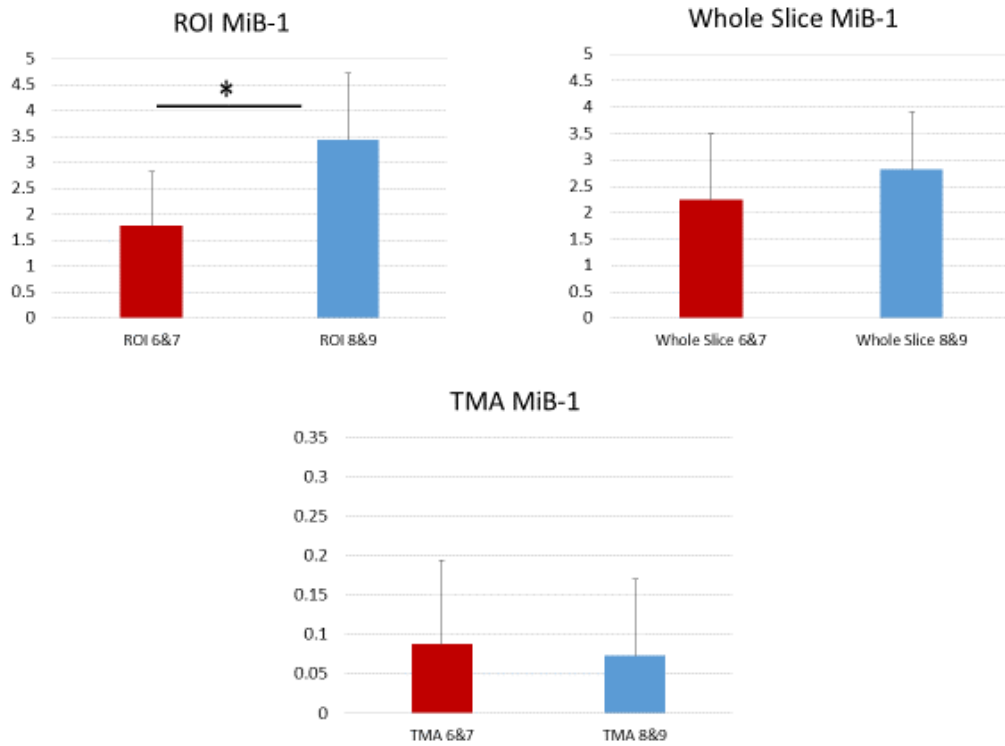


FIGURE 17 – MiB-1 Results

TABLE I

MiB-1 RESULTS

MiB-1			
	6&7 Average (StDev)	8&9 Average (StDev)	Statistically Significant?
Regions of Interest	1.7877 (1.0339)	3.4407 (1.2902)	Yes
Whole Slices	2.2392 (1.2554)	2.8164 (1.0900)	No
Tissue Micro Array	0.0876 (0.1062)	0.0738 (0.0970)	No

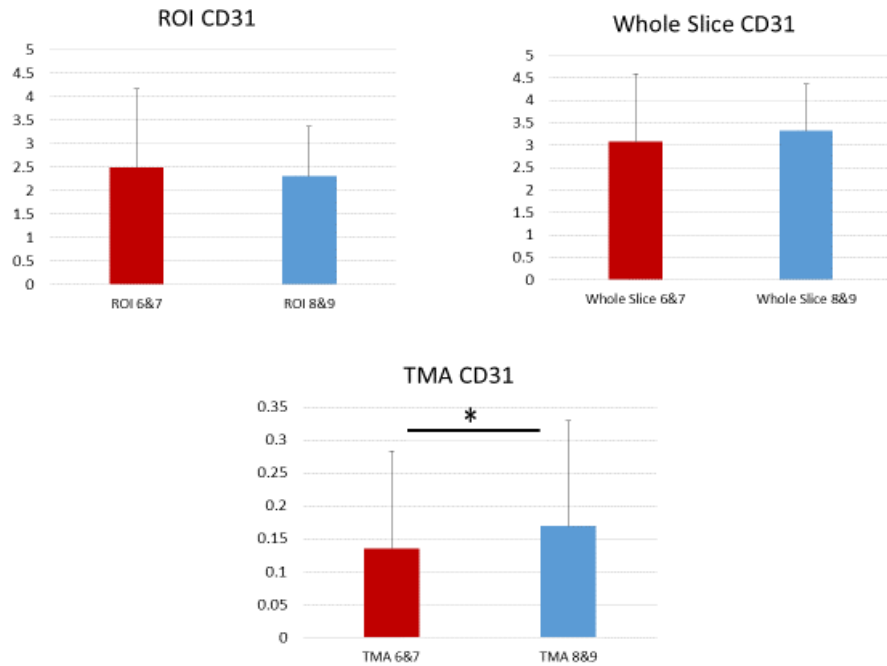


FIGURE 18 – CD31 Results

TABLE II

CD31 RESULTS

CD31			
	6&7 Average (StDev)	8&9 Average (StDev)	Statistically Significant?
Regions of Interest	2.5050 (1.6626)	2.3213 (1.0534)	No
Whole Slices	3.0808 (1.5060)	3.3273 (1.0477)	No
Tissue Micro Array	0.1354 (0.1472)	0.1691 (0.1608)	Yes

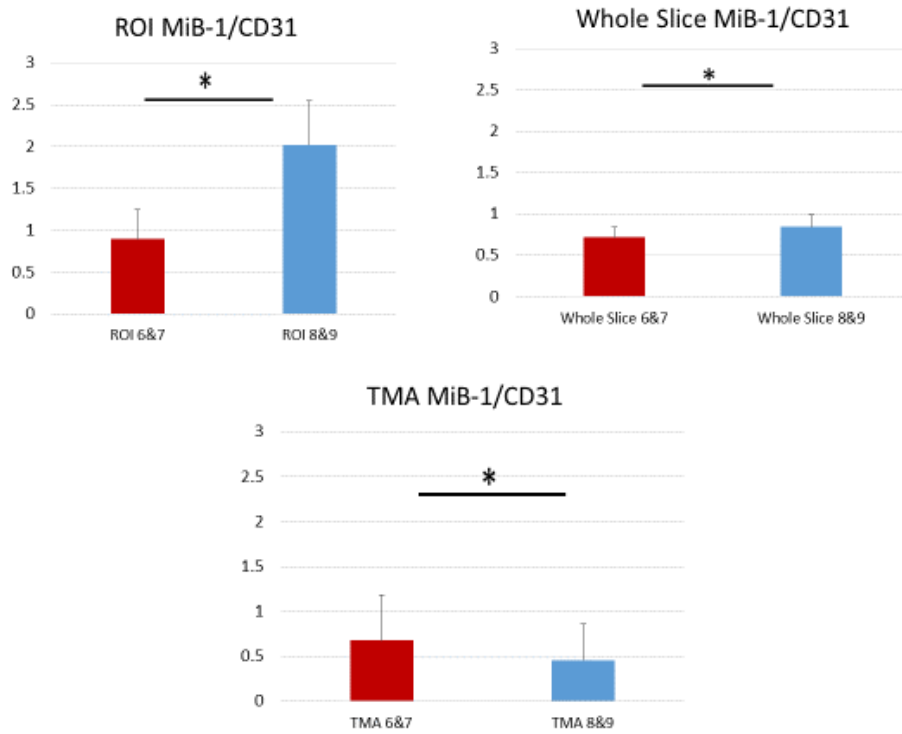


FIGURE 19 – MiB-1/CD31 Ratio Results

TABLE III

MiB-1/CD31 RATIO RESULTS

MiB-1/CD31			
	6&7 Average (StDev)	8&9 Average (StDev)	Statistically Significant?
Regions of Interest	0.8918 (0.3543)	2.0209 (0.5251)	Yes
Whole Slices	0.7161 (0.1270)	0.8451 (0.1410)	Yes
Tissue Micro Array	0.6795 (0.4999)	0.4559 (0.4059)	Yes

The calculations of the t-tests were done using the formula:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{S_1^2}{N_1} + \frac{S_2^2}{N_2}}}$$

where t is the t-score, x represents the averages of the category, S the standard deviation of the category, and N the number of samples in the category, with the subscript 1 categorizing the data with Gleason Scores of 8 or 9 and the subscript 2 categorizing the data with Gleason Scores of 6 or 7.

The degrees of freedom for the t-test was determined by the formula:

$$DoF = N_1 + N_2 - 2$$

The results of the t-test were evaluated at the confidence level of 0.05. The complete tables of the calculation process can be seen below.

TABLE IV

ROI STATISTICAL SIGNIFICANCE

Regions of Interest				
<b>MiB-1/CD31</b>	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>
6/7 AVG	13	0.892	0.354	0.125
8/9 AVG	9	2.021	0.525	0.276
		<b>Variance of Difference</b>	<b>Standard Deviation</b>	
		0.040	0.201	
		<b>t</b>	<b>DoF</b>	<b>p</b>
		5.624739336	20	0.000
SIGNIFICANT				
<b>MiB-1</b>	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>
6/7 AVG	13	1.788	1.034	1.069
8/9 AVG	9	3.441	1.290	1.665
		<b>Variance of Difference</b>	<b>Standard Deviation</b>	
		0.267	0.517	
		<b>t</b>	<b>DoF</b>	<b>p</b>
		3.198	20	0.002
SIGNIFICANT				
<b>CD31</b>	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>
6/7 AVG	13	2.505	1.663	2.764
8/9 AVG	9	2.321	1.053	1.110
		<b>Variance of Difference</b>	<b>Standard Deviation</b>	
		0.336	0.580	
		<b>t</b>	<b>DoF</b>	<b>p</b>
		-0.317	20	0.488
NOT SIGNIFICANT				

TABLE V

WHOLE SLICE STATISTICAL SIGNIFICANCE

<b>Whole Slice Samples</b>				
<b>MiB-1/CD31</b>	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>
6/7 AVG	9	0.716	0.127	0.016
8/9 AVG	8	0.845	0.141	0.020
<b>Variance of Difference</b>			<b>Standard Deviation</b>	
0.004			0.065	
<b>t</b>	<b>DoF</b>	<b>p</b>		
	1.972	15	0.034	SIGNIFICANT
<b>MiB-1</b>	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>
6/7 AVG	9	2.239	1.255	1.576
8/9 AVG	8	2.816	1.090	1.188
<b>Variance of Difference</b>			<b>Standard Deviation</b>	
0.324			0.569	
<b>t</b>	<b>DoF</b>	<b>p</b>		
	1.015	15	0.163	NOT SIGNIFICANT
<b>CD31</b>	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>
6/7 AVG	9	3.081	1.506	2.268
8/9 AVG	8	3.327	1.048	1.098
<b>Variance of Difference</b>			<b>Standard Deviation</b>	
0.389			0.624	
<b>t</b>	<b>DoF</b>	<b>p</b>		
	0.395	15	0.349	NOT SIGNIFICANT



TABLE VI

TMA STATISTICAL SIGNIFICANCE

<b>Tissue Micro Array</b>					
<b>MiB-1</b>					
<b>/CD31</b>	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>	
6/7 AVG	975	0.679	0.500	0.250	
8/9 AVG	132	0.456	0.406	0.165	
		<b>Variance of Difference</b>		<b>Standard Deviation</b>	
		0.002		0.039	
		<b>t</b>	<b>DoF</b>	<b>p</b>	
		5.763	1105	0.000	
				SIGNIFICANT	
<b>MiB-1</b>					
	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>	
6/7 AVG	1213	0.088	0.106	0.011	
8/9 AVG	158	0.074	0.097	0.009	
		<b>Variance of Difference</b>		<b>Standard Deviation</b>	
		6.883E-05		8.296E-03	
		<b>t</b>	<b>DoF</b>	<b>p</b>	
		1.666	1369	0.096	
				NOT SIGNIFICANT	
<b>CD31</b>					
	<b>n</b>	<b>Average</b>	<b>Standard Deviation</b>	<b>Variance</b>	
6/7 AVG	1142	0.135	0.147	0.022	
8/9 AVG	156	0.169	0.161	0.026	
		<b>Variance of Difference</b>		<b>Standard Deviation</b>	
		0.00018463		0.013587876	
		<b>t</b>	<b>DoF</b>	<b>p</b>	
		2.481	1296	0.013	
				SIGNIFICANT	

## V. DISCUSSION

Visual inspection of the whole slices of the prostate show that there tended to be a high degree of heterogeneity within some of the stained samples of higher Gleason scores (i.e. 8 & 9), meaning that the stain was found to be densely populated in certain regions within the slice while being almost non-existent in other regions. It is believed that this is a contributing factor in the results found for the tissue micro-array samples. While the region of interest samples were selected based on the opinion that they were areas representative of the whole tissue, the TMA samples seemed to have a large amount of samples that possess next to no staining, accounting for the lower average stain percentages and the relatively large standard deviations in comparison to the other selection methods.

Taking this into consideration, it is believed that the results seen in the whole slice analysis and the analysis of the regions of interest should be the point of focus moving

forward. The results of these sample types show that samples coming from patients rated at higher Gleason scores (8's and 9's) tend to express higher levels of the MiB-1 stain relative to the CD31 stain in comparison to samples coming from patients rated at lower Gleason scores (6's and 7's).

## VI. CONCLUSION

It has been discussed that a limitation of tissue microarrays is that, due to the small size of the samples taken, the tissue selected for analysis may not be representative of the tissue as a whole. Additionally, this is thought to be a larger problem for epithelial tumors due to the heterogeneous nature of the tumors (7). Prostate cancer primarily takes the form of adenocarcinomas which are formed in the epithelial tissue in the prostate glands which creates a cause for concern for the tissue microarrays used in this study may not have been representative of the larger tissue used to calculate the Gleason score for the sample. It is for this reason that it is believed that primary conclusions from this study should be drawn from the results of the whole tissue and region of interest analysis.

It is the intended use of this program to be used by doctors independently as a supplemental means for helping them better understand the nature of their patients'

condition. With more knowledge regarding the status of the patients' condition, hopefully doctors will be able to make better informed decisions as to what is the best way to treat their patients.

## VII. RECOMMENDATIONS

One of the remaining limitations of this program is its inability to detect imperfections in the source image. If there is a problem that occurred in the generation of the image (e.g. sample folded over, smudges on microscope lens, foreign object in field of view, etc...) the program currently has no way of flagging these samples and can report back inaccurate data. This is a simple process in the manual method due to the ability of the person performing the analysis being able to visually see what's going on before the analysis and can either account for the issue or preemptively remove the sample before it gets put in with the other data. It is therefore recommended that future work on this program be involved with ensuring that there is a way for the program to flag samples that could possess imperfections so that the user can then decide whether those samples are included in their final data or not.

It is believed that the program developed for this study could be adapted for alternate future applications, especially for the use in analysis of other types of tumors

besides prostate tumors and other types of immunohistochemical stains. The primary adjustments that would need to be made to adjust the program for these situations would involve changing the inclusion/exclusion criteria for what constitutes stain vs normal tissue in the slides. The process by which the program goes through the various images and analyzes them should remain the same regardless of the content of the images. Some adjustment to how the analyzed data is handled may also be necessary depending on the desired information.

Additionally, if the program was to be used on a more powerful computer, the program could be optimized further by eliminating/changing certain steps in the process that were necessary due to the limitations imposed by the computer the program was designed on.

## REFERENCES CITED

1. "What Is Prostate Cancer?" *What Is Prostate Cancer?* N.p., n.d. Web. 24 Mar. 2015.
2. "Key Statistics for Prostate Cancer." *What Are the Key Statistics about Prostate Cancer?* N.p., n.d. Web. 24 Mar. 2015
3. "Prostate Cancer Risk Factors." *What Are the Risk Factors For Prostate Cancer?* N.p., n.d. Web. 24 Mar. 2015
4. "Prostate Cancer Treatment." *How is Prostate Cancer Treated?* N.p., n.d. Web. 24 Mar. 2015
5. "Special Stains & Immunohistochemistry." CCPathology. Central Coast Pathology Laboratory, n.d. Web. 23 Feb. 2016
6. "Immunohistochemistry (IHC) in Cancer." *Immunohistochemistry.* N.p., n.d. Web. 10 Mar. 2016
7. KHOUJA, M. HAYSAM, Mark Baekelandt, Agkha Sarab, JAHN M. NESLAND, and Ruth Holm. "Limitations of Tissue Microarrays Compared with Whole Tissue Sections in Survival Analysis." *Oncology Letters.* D.A. Spandidos, 01 Sept. 2010. Web. 22 Mar. 2016.
8. Rizzardi, Anthony E., Arthur T. Johnson, Rachel Isaksson Vogel, Stefan E. Pambuccian, Jonathan Henriksen, Amy PN Skubitz, Gregory J. Metzger, and Stephen C. Schmechel. "Quantitative Comparison of Immunohistochemical Staining Measured by Digital Image Analysis versus Pathologist Visual Scoring." *Diagnostic Pathology.* BioMed Central, 20 June 2012. Web. 07 Apr. 2016.



9. Lan, Chunyan, Andreas Heindl, Xin Huang, Shaoyan Xi, Susana Banerjee, Jihong Liu, and Yinyin Yuan. "Quantitative Histology Analysis of the Ovarian Tumour Microenvironment." *Nature*. Nature, 17 Nov. 2015. Web. 12 Apr. 2016.
10. "CRImage - Tumour Image Analysis." MarkowitzLab. N.p., n.d. Web. 12 Apr. 2016.
11. Alyassin, Mohamad A., SangJun Moon, Hasan O. Keles, Fahim Manzur, Richard L. Lin, Edward Hæggestrom, Daniel R. Kuritzkes, and Utkan Demirci. "Rapid Automated Cell Quantification on HIV Microfluidic Devices." *Lab on a Chip*. U.S. National Library of Medicine, 30 Sept. 2009. Web. 26 Apr. 2016.
12. Bredies, Kristian, Marcus Wagner, Christian Schubert, and Peter Ahnelt. "Computer-assisted Counting of Retinal Cells by Automatic Segmentation after TV Denoising." *BMC Ophthalmology*. BioMed Central, 20 Oct. 2013. Web. 26 Apr. 2016.
13. Al-Khazraji, Baraa K., Philip J. Medeiros, Nicole M. Novielli, and Dwayne N. Jackson. "An Automated Cell-counting Algorithm for Fluorescently-stained Cells in Migration Assays." *Biological Procedures Online*. BioMed Central, 19 Oct. 2011. Web. 26 Apr. 2016.