

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

Article Type: Research Article

Urinary Bladder Cancer Staging in CT Urography using Machine Learning

Sankeerth S. Garapati
Lubomir Hadjiiski
Kenny H. Cha
Heang-Ping Chan
Elaine M. Caoili
Richard H. Cohan
Alon Weizer*
Ajjai Alva**
Chintana Paramagul
Jun Wei
Chuan Zhou

Department of Radiology, The University of Michigan, Ann Arbor, Michigan 48109

**Department of Urology, Comprehensive Cancer Center, The University of Michigan, Ann Arbor, Michigan 48109*

***Department of Internal Medicine, Hematology-Oncology, The University of Michigan, Ann Arbor, Michigan 48109*

Running Title: Bladder Cancer Staging by Machine Learning

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/mp.12510](https://doi.org/10.1002/mp.12510)

This article is protected by copyright. All rights reserved

29 Correspondence:

30

31 Lubomir Hadjiiski

32 Department of Radiology

33 University of Michigan

34 1500 E. Medical Center Drive

35 MIB C476

36 Ann Arbor, MI 48109-5842

37 Telephone: (734) 647-7428

38 Fax: (734) 615-5513

39 E-mail:lhadjisk@umich.edu

40

41 **ABSTRACT**

42 **Purpose:** To evaluate the feasibility of using an objective computer aided system to assess
43 bladder cancer stage in CT Urography (CTU).

44 **Materials and Methods:** A data set consisting of 84 bladder cancer lesions from 76 CTU cases
45 was used to develop the computerized system for bladder cancer staging based on machine
46 learning approaches. The cases were grouped into two classes based on pathological stage $\geq T2$
47 or below T2, which is the decision threshold for neoadjuvant chemotherapy treatment clinically.
48 There were 43 cancers below stage T2 and 41 cancers at stage T2 or above. All 84 lesions were
49 automatically segmented using our previously developed auto-initialized cascaded level sets (AI-
50 CALS) method. Morphological and texture features were extracted. The features were divided
51 into subspaces of morphological features only, texture features only, and a combined set of both
52 morphological and texture features. The data set was split into Set 1 and Set 2 for two-fold cross
53 validation. Stepwise feature selection was used to select the most effective features. A linear
54 discriminant analysis (LDA), a neural network (NN), a support vector machine (SVM), and a
55 random forest (RAF) classifier were used to combine the features into a single score. The
56 classification accuracy of the four classifiers was compared using the area under the receiver
57 operating characteristic (ROC) curve (A_z).

58 **Results:** Based on the texture features only, the LDA classifier achieved a test A_z of 0.91 on Set
59 1 and a test A_z of 0.88 on Set 2. The test A_z of the NN classifier for Set 1 and Set 2 were 0.89

60 and 0.92, respectively. The SVM classifier achieved test A_z of 0.91 on Set 1 and test A_z of 0.89
61 on Set 2. The test A_z of the RAF classifier for Set 1 and Set 2 was 0.89 and 0.97, respectively.
62 The morphological features alone, the texture features alone, and the combined feature set
63 achieved comparable classification performance.

64 **Conclusion:** The predictive model developed in this study shows promise as a classification tool
65 for stratifying bladder cancer into two staging categories: greater than or equal to stage T2 and
66 below stage T2.

67
68 **Keywords:** Radiomics, Computer-Aided Diagnosis, CT Urography, Bladder Cancer Staging,
69 Segmentation, Feature Extraction, Classification, Machine Learning.

72 1. INTRODUCTION

73 Bladder cancer is one of the most common cancers affecting both men and women¹. It can cause
74 substantial morbidity and mortality among the patients with the disease. In 2017, it is estimated
75 that there will be 79,030 new cases and 16,870 deaths from bladder cancer¹. One in 42
76 Americans will be diagnosed with bladder cancer in their lifetime and 9 out of 10 patients with
77 this cancer are over the age of 55^{1,2}. The average age of diagnosis is 73¹. Approximately half of
78 all bladder cancer cases are first found while the cancer is still confined to the inner wall of the
79 bladder and has not invaded into deeper layers or distant parts of the body¹. Bladder cancer has a
80 recurrence rate of 50-80 percent and requires constant surveillance. This makes it the most
81 expensive cancer to treat, requiring a total of \$4.1 billion yearly, on a per patient basis in the
82 United States². Bladder cancer can be divided into three categories that include noninvasive,
83 superficial, and invasive. The initial treatment for bladder cancer is transurethral resection of the
84 bladder tumor (TURBT), which removes the tumor from the bladder and also helps provide
85 information regarding the stage of the cancer³⁻⁵. Bladder cancer is staged in order to determine
86 treatment options and estimate a prognosis for the patient. Accurate staging provides the
87 physician with information about the extent of the cancer. The tumor stages T refer to the depth
88 of the penetration of the tumor into the layers of the bladder. T0 indicates no primary tumor, T1
89 indicates that the tumor has invaded the connective tissue under the epithelium, T2 indicates that
90 the tumor has invaded the bladder muscle, T3 indicates that the tumor has invaded the fatty

91 tissue around the bladder, and T4 indicates that the tumor has spread beyond the fatty tissue into
92 other areas such as the pelvic wall, uterus, prostate or abdominal wall⁶ (Fig. 1). An example of
93 bladder cancer stage T2 is presented in Fig. 2.

94 The accurate staging of bladder cancer is crucial to providing proper treatment to the patient.
95 Superficial diseases (under stage T2) can be managed with less aggressive treatment than
96 invasive diseases (stage T2 and above)³⁻⁵. There are two types of staging for bladder cancer -
97 clinical and pathological. The clinical stage is the physicians' best estimate for the extent of the
98 cancer based on physical exams and imaging. The pathological stage is determined by analysis
99 of the tissue collected from the cancer after biopsy, tumor resection or bladder cystectomy. The
100 accuracy of the staging depends on the complete resection of the tumor. Incomplete resection of
101 the tumor may reduce the reliability of the staging at the beginning of the tumor management
102 process⁷. Bladder cystectomy ensures that the entire bladder tumor is present for pathological
103 review; therefore, the pathological staging is based on the histological review of the cystectomy
104 specimen⁶. Adjuvant chemotherapy is used in patients with locally advanced bladder cancer in
105 order to reduce the chances of cancer recurrence following radical cystectomy⁸. Neoadjuvant
106 chemotherapy is used prior to radical cystectomy in order to reduce the tumor size before
107 surgical removal; for example, a cisplatin-based regimen has been shown to decrease the
108 probability of finding extravesical disease and improve survival when compared to radical
109 cystectomy alone⁸⁻¹⁰.

Author Manuscript

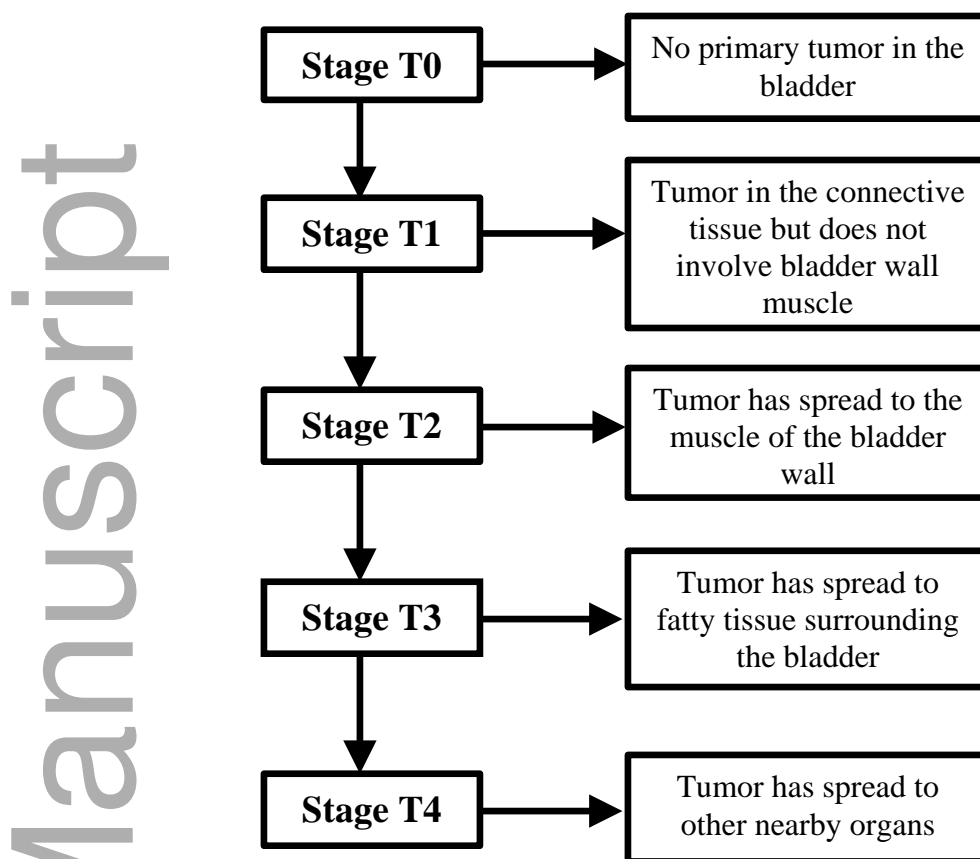


Figure 1. Bladder cancer stage grading scale definition.

110
 111 Correct staging of bladder cancer is crucial for the decision of neoadjuvant chemotherapy
 112 treatment and minimizing the risk of under-treatment or over-treatment. Patients with stage T2 to
 113 T4 carcinomas of the bladder are recommended for treatment with neoadjuvant chemotherapy.
 114 Studies found that up to 50% of the patients who are estimated to have a T1 disease at clinical
 115 staging are under-staged and later upstaged after radical cystectomy¹¹⁻¹⁴. This inaccuracy in
 116 staging can partly be attributed to the subjectivity and variability of clinicians in utilizing various
 117 diagnostic information. The purpose of this study is to develop an objective decision support
 118 system that can potentially reduce the risk of under-treatment or over-treatment by merging
 119 radiomic information in a predictive model using statistical outcomes and machine learning.
 120

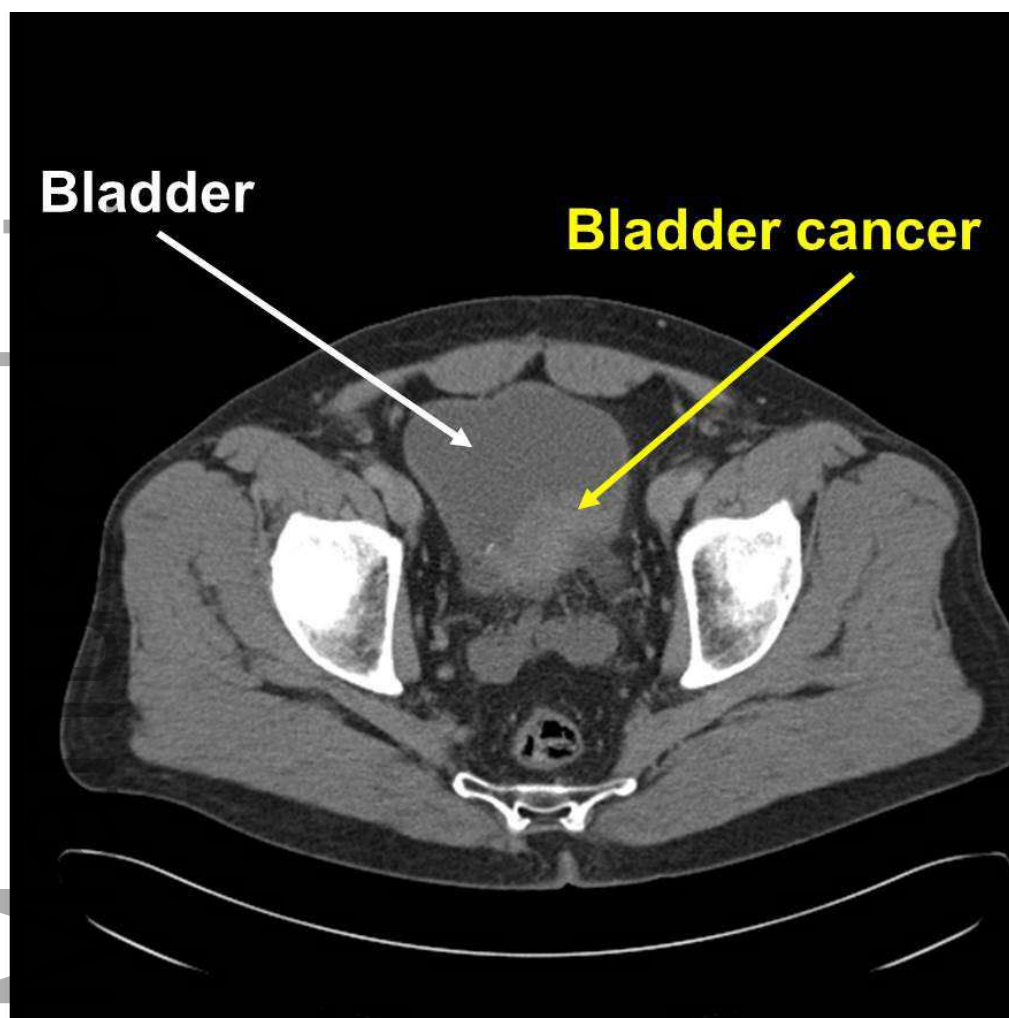


Figure 2. Urinary Bladder CT. The bladder cancer is marked and clearly visible. The cancer stage is T2.

121

122

123

124

2. MATERIALS AND METHODS

2.1 Data Set

The data collection protocol was approved by our institutional review board and is HIPAA compliant. Patient informed consent was waived for this retrospective study. Our data set consisted of 84 bladder cancer lesions from 76 bladder cancer CTU cases collected from patient files without additional imaging for research purpose. The CTU scans in this data set were acquired at an image slice interval of 0.625 to 1.25 mm using 120 kVp and 120-280 mA. The

130

131 data set consisted of 22 non-contrast cases (22 lesions), 22 early phase contrast-enhanced cases
132 (22 lesions), and 32 delayed-phase contrast-enhanced cases (40 lesions). Per imaging protocol,
133 the early phase contrast-enhanced images are obtained 60 seconds following the initiation of a
134 contrast injection. The delayed-phase contrast-enhanced images are obtained 12 min after the
135 initiation of contrast injection. The type of scan a patient receives is determined by the protocol
136 of the hospital performing the scan. Our data set includes patients referred to our hospital for
137 treatment so that some scans were performed at outside hospitals and followed different scanning
138 protocols, resulting in scans with inconsistent contrast-enhancement phase. A patient may also
139 get a non-contrast scan due to risk factors, such as allergy to the contrast media, asthma, renal
140 insufficiency, significant cardiac disease, or anxiety¹⁵.

141 For all cases, clinical and pathological staging were performed during the patient's
142 clinical care. Cystectomy was performed after completing the course of neoadjuvant
143 chemotherapy. The primary chemotherapy regimen used for the patients in our data set were
144 MVAC, which is a combination of four medications: Methotrexate, Vinblastine, Doxorubicin,
145 and Cisplatin. Stage T2 is identified to be clinically important as a decision threshold for
146 neoadjuvant chemotherapy treatment. The stage at the beginning of the tumor management
147 process, based on the clinical staging and pathological staging was used as a reference standard
148 of the tumor stage for our study.

149 In addition, for all bladder cancer lesions a radiologist measured the longest diameter on
150 the pre-treatment scans by using an electronic caliper provided by an in-house developed
151 graphical user interface.

152 The 84 bladder cancer lesions were separated into two classes. The first class consisted
153 of 41 cancers that were stage T2 or above and the patients were treated with neoadjuvant
154 chemotherapy. The second class consisted of 43 cancers that were below stage T2 and patients
155 were not referred to neoadjuvant chemotherapy treatment. The data set was then split randomly
156 by case into two sets with 42 cancers each while keeping the proportion of cancers between the
157 two classes similar. The first set (Set 1) consisted of 22 cancers below stage T2 and 20 cancers
158 stage T2 or above. The second set (Set 2) consisted of 21 cancers below stage T2 and 21 cancers
159 stage T2 or above.

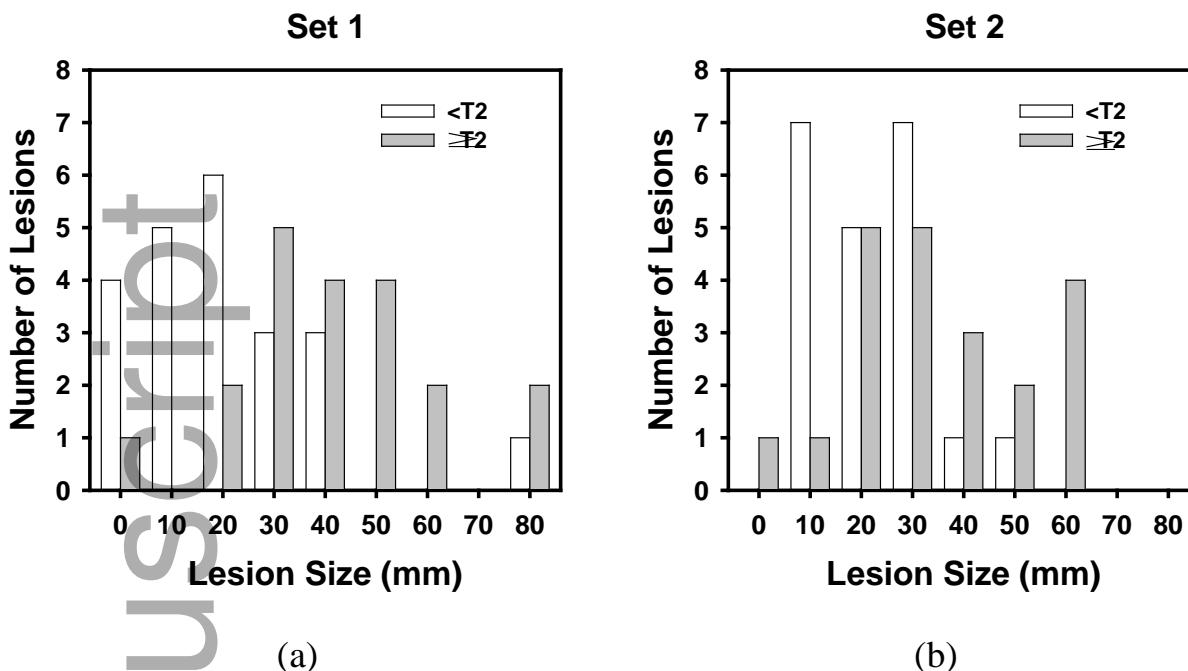


Figure 3. Distribution of tumor sizes (the longest diameters) for Set 1 and Set 2. (a) Set 1: The average tumor sizes of stage < T2 and \geq T2 were 26.4 ± 17.3 mm and 45.6 ± 19.1 mm respectively. (b) Set 2: The average tumor sizes of stage < T2 and \geq T2 were 27.3 ± 10.8 mm and 40.6 ± 17.3 mm respectively.

160

161 In Set 1, two patients had two lesions and one patient had three lesions. In Set 2, three patients
 162 had two lesions. In Set 1, the average tumor sizes (the longest diameters) of stage <T2 and \geq T2
 163 were 26.4 ± 17.3 and 45.6 ± 19.1 mm, respectively (Fig. 3a). In Set 2, the average tumor sizes (the
 164 longest diameters) of stage <T2 and \geq T2 were 27.3 ± 10.8 mm and 40.6 ± 17.3 mm, respectively
 165 (Fig. 3b).

166

167 2.2 Segmentation of Bladder Lesions on CT Urography

168 Our previously developed method for bladder lesion segmentation using an auto-initialized
 169 cascaded level set (AI-CALS) was used¹⁶. Briefly, the system consists of three stages that include
 170 preprocessing, initial segmentation, and 3D level set segmentation (Fig. 4). The segmentation of
 171 bladder lesions is often difficult as some lesions are located in the non-contrast enhanced region
 172 of the bladder such that contrast between the lesion and the surrounding background was low.
 173 Additionally, lesions often have irregular boundaries and can be very small and subtle. Each

174 lesion in the data set was marked by a bounding box as an input volume of interest (VOI). The
175 lateral dimensions of the box were determined by an adjustable rectangle within the image slice
176 that contains the best view of the lesion. The top and bottom slices are marked to completely
177 enclose the lesion. The AI-CALS segmentation is then automatically performed in the VOI. In
178 the pre-processing stage, image processing techniques including smoothing, anisotropic
179 diffusion, gradient filters, and a rank transform of the gradient magnitude are used to generate
180 sets of smoothed images, gradient magnitude images, and gradient vector images. The initial
181 segmentation surface is obtained by combining information from these images. Three
182 dimensional (3D) flood fill algorithm, morphological dilation filter, and morphologic erosion
183 filter are applied to the initial segmentation surface to connect nearby components, which is then
184 used to initialize the level set segmentation. The initial contour is propagated toward the lesion
185 boundary using a bank of cascaded level sets. The level sets help refine the initial contour. The
186 details of the AI-CALS method can be found in our previous paper¹⁶.
187

Author Manuscript

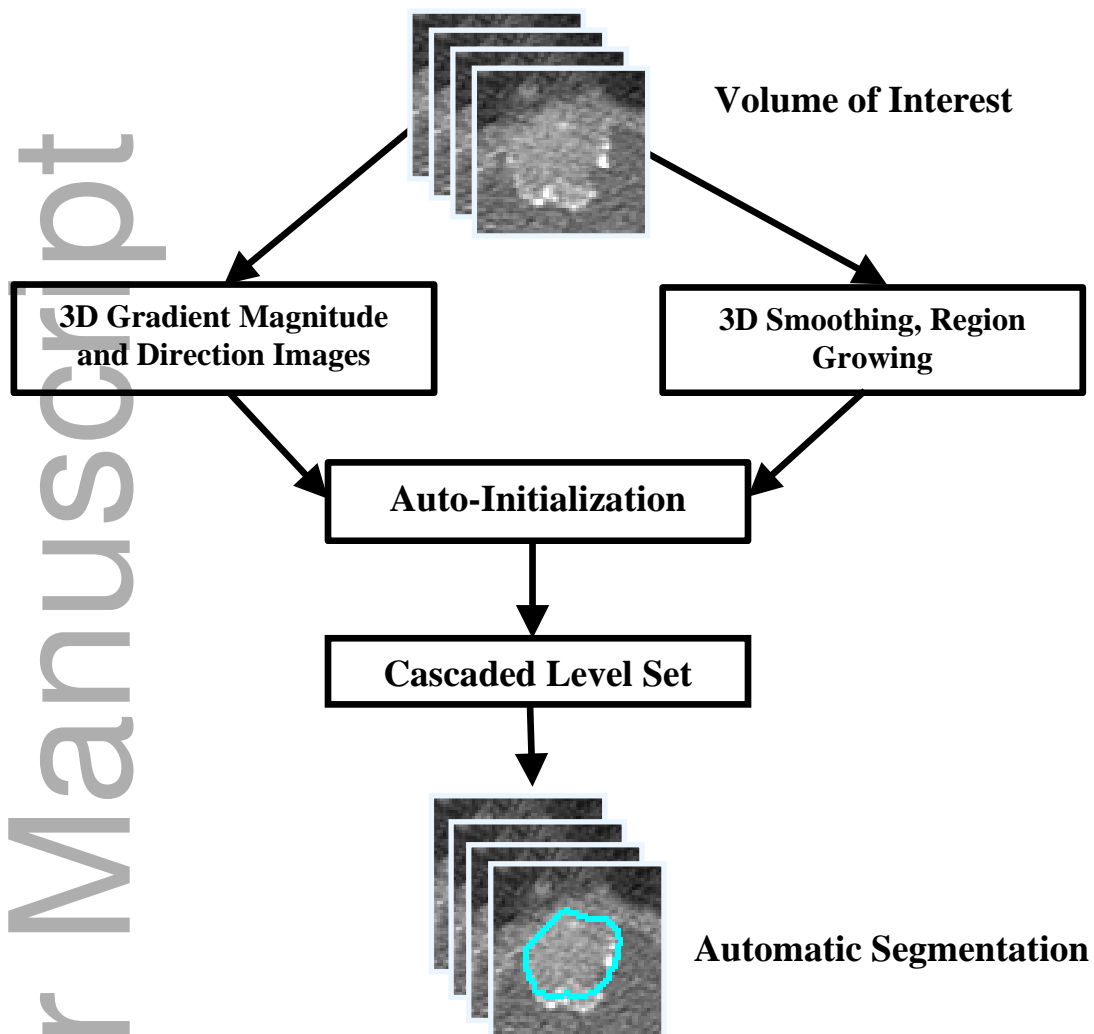


Figure 4. Block diagram of the auto-initialized cascaded level sets (AI-CALS) method.

188

189

190

191

3. CLASSIFICATION

3.1 Feature Extraction

193

Following automated computer segmentation, texture features and morphological

194

features were extracted to characterize the lesion. The mass size was measured as its 3D volume.

195

Five morphological features were extracted based on the normalized radial length (NRL). NRL is

196 defined as the radial length normalized relative to the maximum radial length for the segmented
197 object¹⁷. The NRL features extracted include zero crossing count, area ratio, standard deviation,
198 mean, and entropy. In addition, ten contrast features and a number of features including
199 circularity, rectangularity, perimeter-to-area ratio, Fourier descriptor, gray level average,
200 standard deviation of gray level, mean density, eccentricity, moment ratio, and axis ratio were
201 extracted as shape descriptors.

202 The texture of the tumor margin can provide important information about its
203 characteristics. We calculated texture features from the rubber band straightening transform
204 (RBST) images¹⁸ of the tumor margin including those from the run-length statistics matrices,
205 filtered Dasarathy east-west direction and filtered Dasarathy horizontal direction^{19,20}. The texture
206 feature set also included the gray level radial gradient direction features.

207 In total, 91 features were extracted to form the feature space, including 26 morphological
208 features and 65 texture features.

209

210 **3.2 Feature Selection/Classification**

211 A block diagram of the machine learning based bladder cancer staging system is shown
212 in Fig. 5. Stepwise feature selection was used to select the best subset of features to create an
213 effective classifier²¹. A number of different classification experiments were performed to
214 determine the best collection of input features. The classification performance was compared in
215 three feature spaces: (1) morphological features only, (2) texture features only, and (3)
216 morphological and texture features combined. A two-fold cross validation was conducted by
217 partitioning the data set into Set 1 and Set 2. In the first fold, Set 1 was used for feature selection
218 and classifier training. The trained classifier was then tested on Set 2. In the second fold, feature
219 selection and classifier training were performed on Set 2 and then tested on Set 1.

220 When training on a given fold (for example, Set 1) a leave-one-case-out resampling
221 scheme with stepwise feature selection was used to reduce the dimensionality of the feature
222 space. In stepwise feature selection, one feature is entered or removed in alternate steps while
223 their effect is analyzed using the Wilks' lambda criterion²¹. The significance of the change in the
224 Wilks' lambda when a feature is included or removed was estimated by F statistics. F_{in} , F_{out} , and
225 tolerance are the parameters of the stepwise feature selection, which define the thresholds for
226 inclusion or exclusion of a given feature. A range of F_{in} , F_{out} , and tolerance values is evaluated

227 by using an automated simplex optimization method. The set of F_{in} , F_{out} , and tolerance values
228 that lead to the highest classification result with the lowest number of features based on the
229 training set are selected. A smaller number of features are preferred in order to reduce the chance
230 of overfitting. Once the set of F_{in} , F_{out} , and tolerance is selected, the stepwise feature selection
231 with the selected parameter set is applied to the entire training fold to select a single set of
232 features and train a single classifier. After the classifier is fixed it is applied to the test fold (for
233 example, Set 2) for performance evaluation.

234 Four different classifiers were evaluated in this study. The same partitioning of Set 1 and
235 Set 2 was used for all classifiers. We compared the four classifiers for this classification task.
236 The first classifier was linear discriminant analysis (LDA)^{22,23}. The LDA with the stepwise
237 feature selection was used to determine the most effective features using the training set in each
238 fold, as described above. The second classifier was a back-propagation neural network (NN)²⁴
239 with a single hidden layer and a single output node. The selected features from LDA were used
240 for this classifier and they determined the number of input nodes to the NN. The parameters for
241 the NN were adjusted using the training set, and the best performing network was applied to the
242 test set. The third classifier was a support vector machine (SVM)^{25,26} with a radial basis kernel.
243 Using training data, a SVM determines a decision hyperplane to separate the two classes by
244 maximizing the distance, or the margin, between the training samples of both classes and the
245 hyperplane. The width of the SVM radial basis kernels γ was varied between 0.02 to 0.14 for the
246 experiments. The best parameters for the SVM kernels for a specific experiment were selected
247 using the training set, which were then applied to the test set. The LDA selected features were
248 also used as the input to the SVM. The fourth one is the Random Forest (RAF) classifier²⁷. We
249 used the WEKA²⁸ implementation and selected 50 to 100 trees and 5 to 7 features per tree for
250 our classification task using the training set in each fold. The parameters for the random forest
251 classifier were determined experimentally using the training sets. All 91 features were used as an
252 input to the RAF.

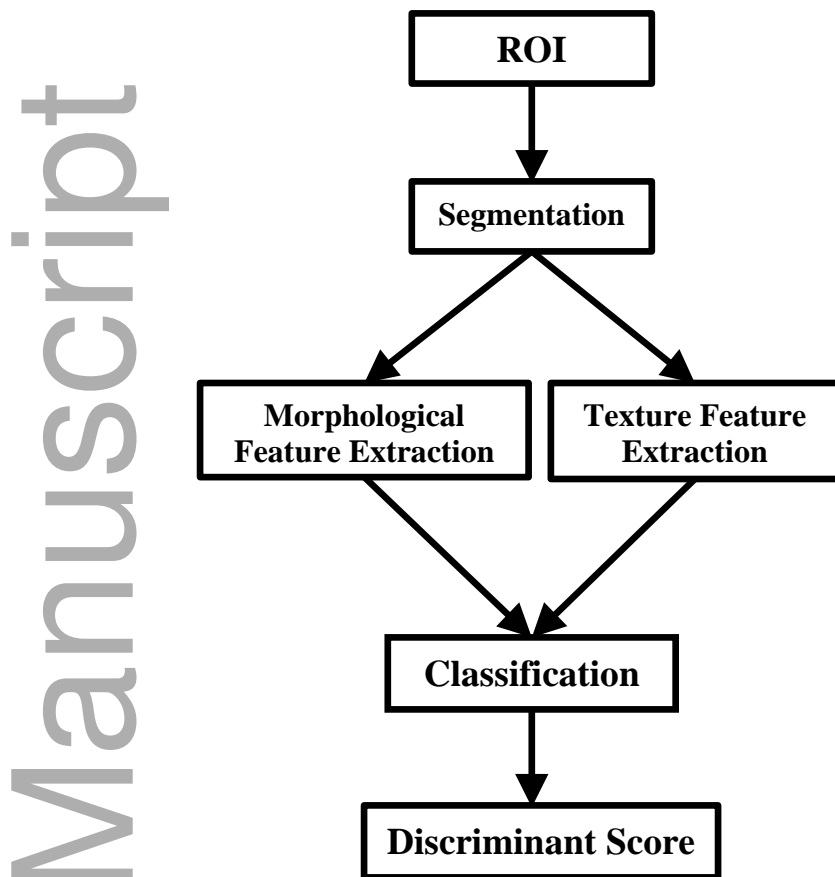


Figure 5. Block diagram of our machine learning based staging system. We compared the linear discriminant analysis (LDA), back-propagation neural network (NN), Support vector machine (SVM), and Random forest classifiers (RAF) in the classification stage for this study.

253

254 3.3 Evaluation Methods

255 Lesion segmentation performance was evaluated using radiologists' 3D hand-segmented
 256 contours as reference standards. The hand outlines of all 84 lesions were obtained from an
 257 experienced abdominal radiologist (RAD1). Hand outlines for a subset of 12 lesions were
 258 obtained from a second experienced abdominal radiologist (RAD2). The average distance and
 259 the Jaccard index²⁹ were calculated between the computer outlines and the hand outlines. The
 260 average distance, *AVDIST*, is defined as the average of the distances between the closest points
 261 of the two contours:

$$AVDIST(G, U) = \frac{1}{2} \left(\frac{\sum_{x \in G} \min\{d(x, y): y \in U\}}{N_G} + \frac{\sum_{y \in U} \min\{d(x, y): x \in G\}}{N_U} \right), \quad (1)$$

262 where G and U are two contours being compared. N_G and N_U denote the number of voxels on G
 263 and U , respectively. The function d is the Euclidean distance. For a given voxel along the
 264 contour G , the minimum distance to a point along the contour U is determined. The minimum
 265 distances obtained for all points along G are averaged. This process is repeated by switching the
 266 roles of G and U . $AVDIST$ is then calculated as the average of the two average minimum
 267 distances.

268 The Jaccard index is defined as the ratio of the intersection between the reference
 269 volume and the segmented volume to the union of the reference volume and the segmented
 270 volume:

$$JACCARD^{3D} = \frac{V_G \cap V_U}{V_G \cup V_U}, \quad (2)$$

271 A value of 1 indicates that V_U completely overlaps with V_G , whereas a value of 0 implies V_U
 272 and V_G are disjoint.

273 To evaluate the classifier performance, the training and test scores output from the
 274 classifier were analyzed using the receiver operating characteristic (ROC) methodology³⁰. The
 275 classification accuracy was evaluated using the area under the ROC curve, A_z . The statistical
 276 significance of the differences between the different classifiers and feature spaces were estimated
 277 by the CLABROC program using ROC software by Metz et al.^{31,32}.

278

279

280

4. RESULTS

281 The lesion segmentation performance of the AI-CALS compared to the radiologist hand outlines
 282 for the 84 lesions are shown in Table 1. Table 2 shows the computer segmentation performance
 283 compared to two different radiologists' hand outlines for a subset of 12 lesions.

Table 1. Segmentation performance of the 84
 lesions compared to hand-outlines performed
 by radiologist 1 (RAD1).

*AI-CALS vs RAD1***Average distance**

4.9 ± 2.7 mm

*AVDIST***Jaccard index**

43.5 ± 14.0%

JACCARD^{3D}

284

Table 2. Segmentation performance for a subset of 12 lesions compared to hand-outlines performed by two different radiologists (RAD1, RAD2)

	<i>AI-CALS vs RAD1</i>	<i>AI-CALS vs RAD2</i>	<i>RAD1 vs RAD2</i>
Average distance			
<i>AVDIST</i>	5.2 ± 2.5 mm	4.1 ± 1.5 mm	2.9 ± 1.1 mm
Jaccard index			
<i>JACCARD</i> ^{3D}	43.2 ± 13.2%	50.1 ± 14.7%	58.7 ± 11.1%

285

286 The performance of the classifiers based on different machine learning techniques, the
 287 LDA, NN, SVM, and RAF, is summarized in Table 3. Different feature spaces containing the
 288 morphological features, the texture features, and the combined set of both morphological and
 289 texture features were used for classification. The features selected with LDA were used in the
 290 SVM and NN classifiers. The LDA classifier with morphological features achieved a training A_z
 291 of 0.91 on Set 1 and a test A_z of 0.81 on Set 2. For training on Set 2 it achieved a A_z of 0.97 and
 292 a test A_z of 0.90 on Set 1. The selected features on the training sets included volume, a contrast
 293 feature, and gray level feature. The test A_z of the NN for Set 1 and Set 2 was 0.88 and 0.91
 294 respectively. The SVM achieved test A_z of 0.88 on Set 1 and test A_z of 0.90 on Set 2. The test
 295 A_z of the RAF for Set 1 and Set 2 was 0.83 and 0.88 respectively. The distribution of the
 296 discriminant scores from the four classifiers for testing on Set 1 and Set 2 in two fold cross-
 297 validation in the morphological feature space are presented in Fig 6. It can be observed that most
 298 of the classifiers were able to provide a relatively good separation between the two classes.

299 By using the texture features the LDA classifier achieved a test A_z of 0.91 on Set 1 and a
 300 test A_z of 0.88 on Set 2. When trained on Set 1 or Set 2 the stepwise feature selection procedure
 301 selected subsets of the filtered Dasarathy east-west direction features, the filtered Dasarathy
 302 horizontal direction features and the gray level radial gradient direction features. The test A_z of
 303 the NN classifier for Set 1 and Set 2 was 0.89 and 0.92, respectively. The SVM classifier
 304 achieved test A_z of 0.91 on Set 1 and test A_z of 0.89 on Set 2. The test A_z of the RAF classifier
 305 for Set 1 and Set 2 was 0.89 and 0.97, respectively.

306 When the morphological and the texture features were combined, the LDA classifier
 307 achieved a test A_z of 0.89 on Set 1 and a test A_z of 0.90 on Set 2. When trained on Set 1 or Set 2
 308 the stepwise feature selection procedure selected a contrast feature, subsets of the filtered
 309 Dasarathy horizontal direction features, and subsets of the gray level radial gradient direction
 310 features. The test A_z of the NN classifier for Set 1 and Set 2 was 0.91 and 0.95, respectively. The
 311 SVM classifier achieved test A_z of 0.92 on Set 1 and test A_z of 0.89 on Set 2. The test A_z of the
 312 RAF classifier for Set 1 and Set 2 was 0.86 and 0.96, respectively. The test ROC curves for all of
 313 the classifiers when tested on Set 1 and Set 2 in the two fold cross-validation in the different
 314 feature spaces are shown in Fig. 7.

315 The differences in the A_z values between pairs of classifiers did not achieve statistical
 316 significance. The classifiers achieved slightly higher A_z values in the texture and combined
 317 feature spaces than in the morphological feature space; however, the differences did not achieve
 318 statistical significance after Bonferroni correction for the multiple comparisons (p -value <
 319 $0.05/18=0.0028$ to be considered significant).

320

321

322 **Table 3.** Summary results for LDA, NN, SVM and RAF classifiers in morphological, texture,
 323 and combined feature spaces. The column “Number of Features” did not apply to the
 324 RAF classifier. All features were used for the RAF classifier. The differences in the A_z
 325 values between pair-wise comparison of the different classifiers did not achieve
 326 statistical significance after performing Bonferroni correction for the 18 comparisons
 327 ($p>0.0028$).

328

		LDA	NN	SVM	RAF
--	--	-----	----	-----	-----

Feature Type	<i>Number of Features</i>	<i>Training</i>	<i>Testing</i>	<i>Training</i>	<i>Testing</i>	<i>Training</i>	<i>Testing</i>	<i>Training</i>	<i>Testing</i>
Morphological Features									
Training (Set 1) Testing (Set 2)	4	0.91	0.81	0.96	0.91	0.95	0.90	1	0.88
Training (Set 2) Testing (Set 1)	4	0.97	0.90	0.98	0.88	0.97	0.88	1	0.83
Texture Features									
Training (Set 1) Testing (Set 2)	2	0.91	0.88	0.95	0.92	0.92	0.89	1	0.97
Training (Set 2) Testing (Set 1)	7	1	0.91	1	0.89	1	0.91	1	0.89
Combined Features									
Training (Set 1) Testing (Set 2)	3	0.92	0.90	0.97	0.95	0.92	0.89	1	0.96
Training (Set 2) Testing (Set 1)	7	1	0.89	1	0.91	1	0.92	1	0.86

329

330

331

332

333

334

335

336

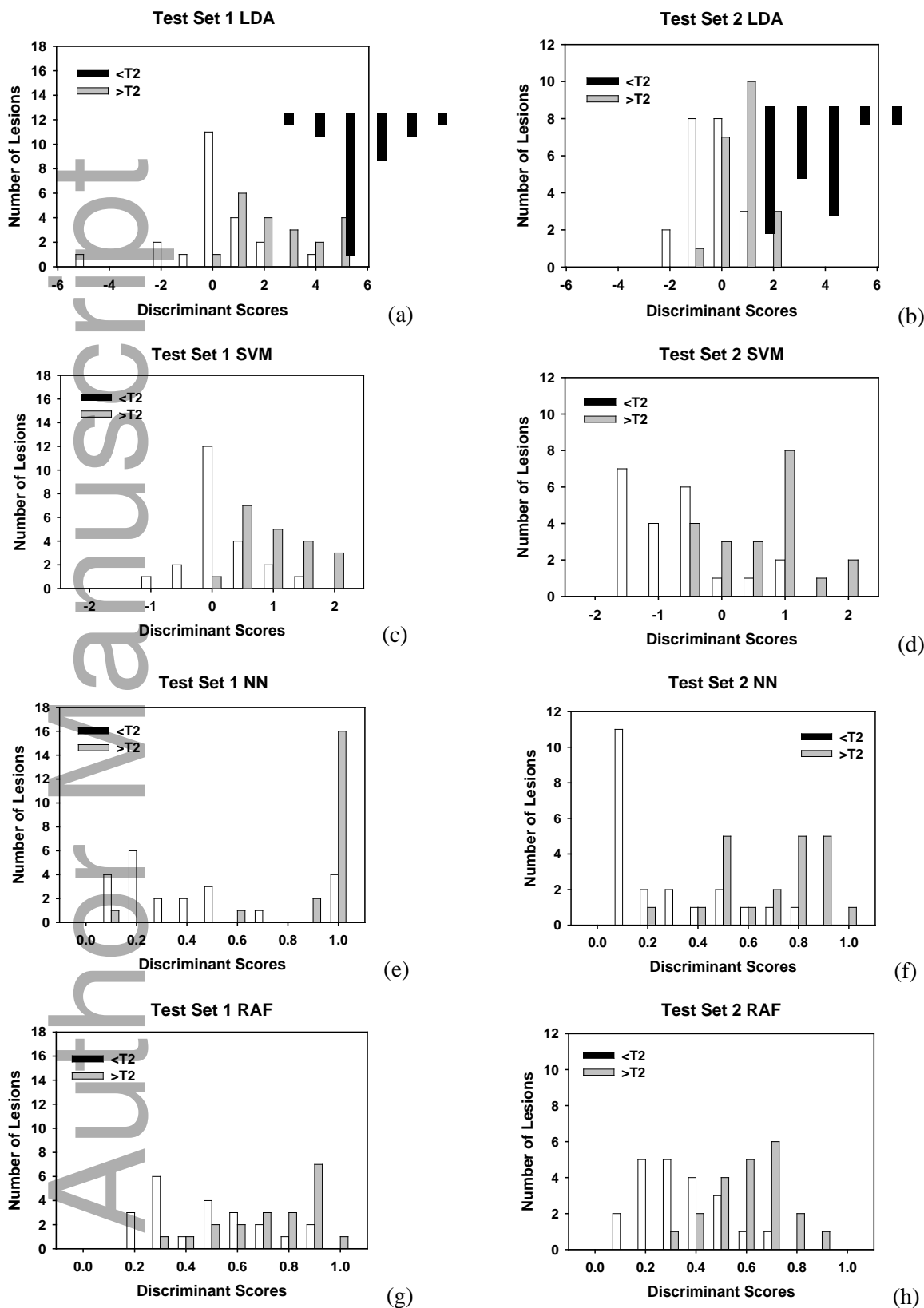


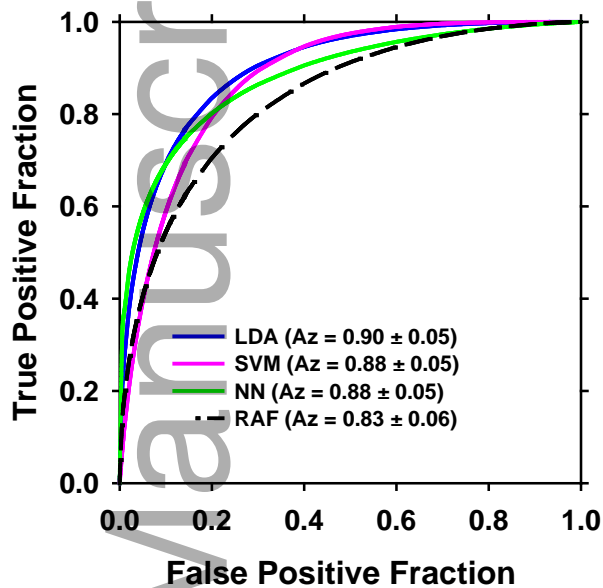
Figure 6. Distribution of the classifiers discriminant scores for testing on Set 1 and Set 2 in two-fold

cross validation using the morphological features. (a) LDA (Set 1) $A_z = 0.90$, (b) LDA (Set 2) $A_z = 0.81$, (c) SVM (Set 1) $A_z = 0.88$, (d) SVM (Set 2) $A_z = 0.90$, (e) NN (Set 1) $A_z = 0.88$, (f) NN (Set 2) $A_z = 0.91$, (g) RAF (Set 1) $A_z = 0.83$, (h) RAF (Set 2) $A_z = 0.88$.

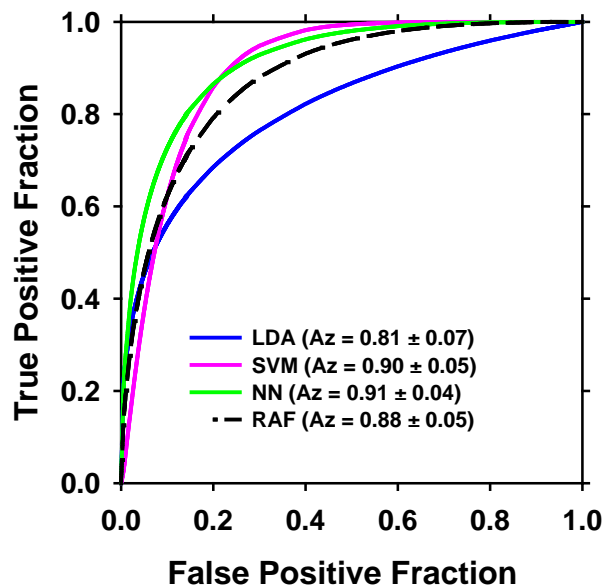
337

338

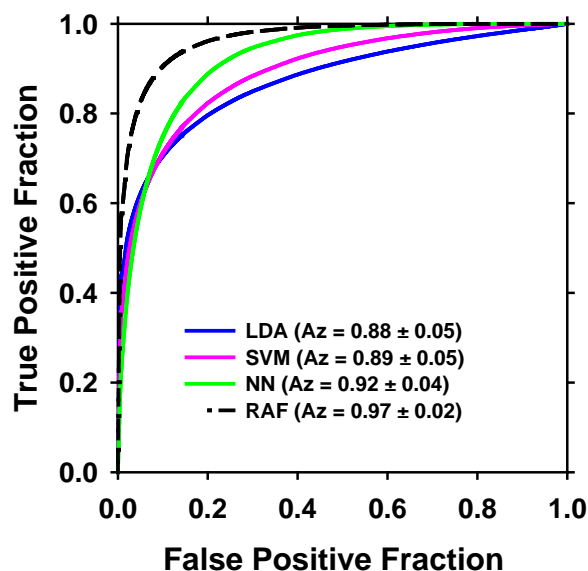
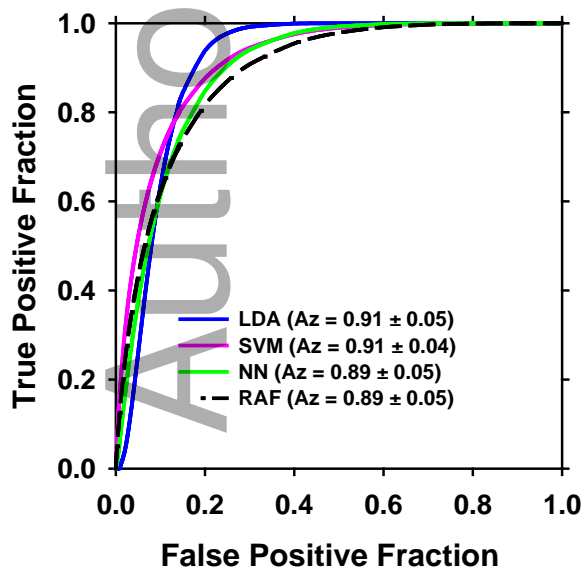
339



(a)



(b)



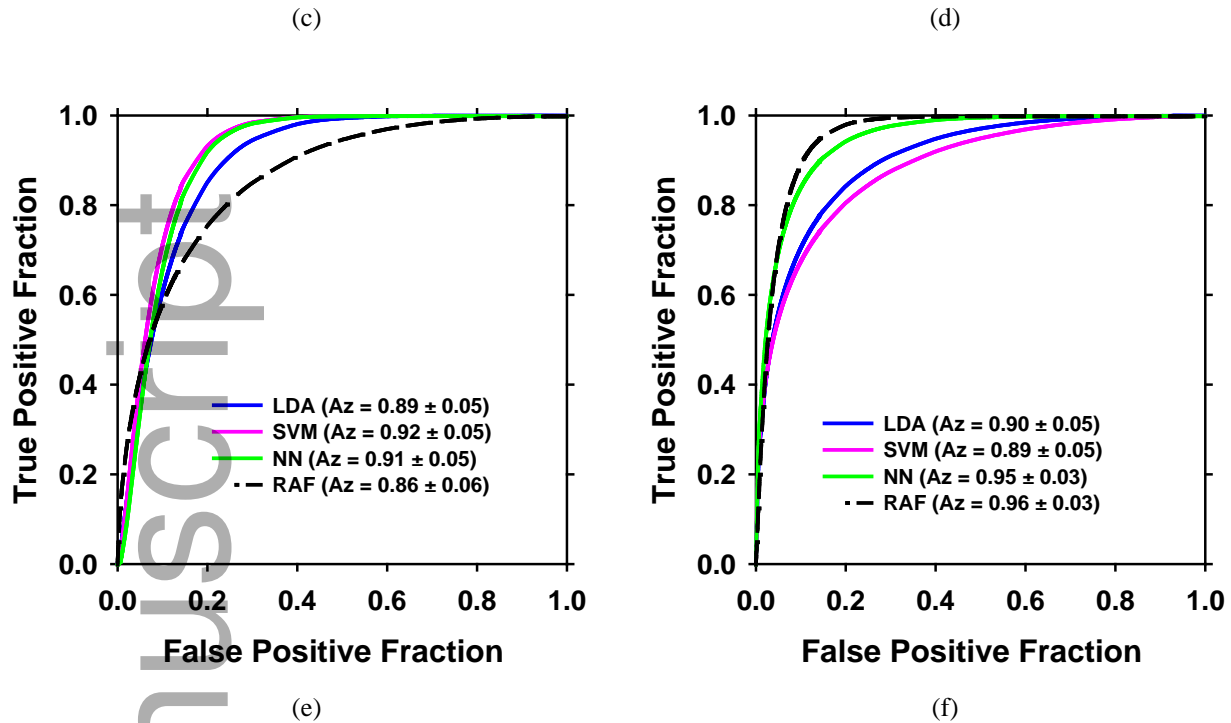


Figure 7. ROC curves for testing on Set 1 and Set 2 in two-fold cross validation for LDA, SVM, NN, and RAF classifiers: Left column: testing on Set 1, right column: testing on Set 2. (a) and (b) morphological features; (c) and (d) texture features; (e) and (f) combined features.

340

341

342

5. DISCUSSION

343

344

345

346

347

348

349

350

351

352

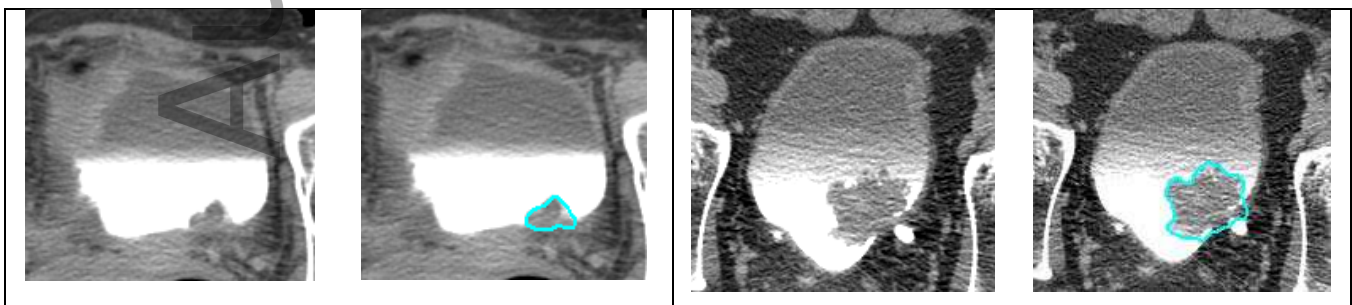
353

The agreement between the AI-CALS lesion segmentation and the radiologists' manual segmentation was slightly lower than the agreement between two radiologists' hand outlines, indicating that the computer segmentation will need to be further improved. Both the morphological and the texture features were important for classifying the bladder cancer stage. When only morphological features were used in the classifier, volume and contrast features were always selected. Volume was the primary feature used to describe lesion size. When the classifier used only the texture features, the features from the 3 main groups, the filtered Dasarathy east-west direction features, the filtered Dasarathy horizontal direction features, and the gray level radial gradient direction features were consistently selected. There was essentially no change in classification accuracy when the morphological features were added to the texture features in the combined set.

354 The LDA, SVM, and NN classifiers all led to relatively consistent results. There was no
 355 statistically significant difference in the performances between pairs of the classifiers. The best
 356 overall results for the two-fold cross validation were obtained when a combined feature set was
 357 used with an NN classifier. Using Set 1 for training, the training A_z was 0.97 and the test A_z was
 358 0.95. Using Set 2 for training, the training A_z was 1.00 and the test A_z was 0.91.

359 The RAF classifier showed greater imbalance between Set 1 and Set 2 than the other
 360 classifiers. When training was done on Set 2 and testing on Set 1, the A_z were substantially
 361 lower than the A_z values when training was done on Set 1 and testing on Set 2. For example, the
 362 test A_z decreased from 0.88 to 0.83 for morphological features, from 0.97 to 0.89 for texture
 363 features only, and from 0.96 to 0.86 for the combined features. This imbalance between the two
 364 sets could be due to the fact that RAF utilized all the features in the subspace whereas the other
 365 three classifiers involved feature selection.

366 Examples of bladder cancers with stages $\geq T2$ or $< T2$ and the corresponding classifier
 367 scores are shown in Fig. 8. The reported scores are test scores for the LDA, SVM, NN, and RAF
 368 classifiers based on the morphological features. In Fig. 8a, b and Fig. 8c, d are shown T1 stage
 369 cancers of different sizes that were correctly classified with low scores by all classifiers. Note
 370 that the output score ranges are different for different classifiers so that the score values should
 371 not be compared across classifiers. T3 stage and T2 stage cancers that were correctly classified
 372 with high scores from all classifiers are presented in Fig. 8e, f and Fig 8g, h, respectively. A case
 373 that was clinically identified as T1 stage pre-surgery but later was identified as a T2 stage cancer
 374 post-surgery is shown in Fig. 8k, l. The classifiers classified the cancer as $\geq T2$ with high scores.
 375 Fig. 8m, n show a T2 stage cancer that was incorrectly identified by the LDA, SVM, and NN
 376 classifiers with low scores, but correctly identified by the RAF with a high score.



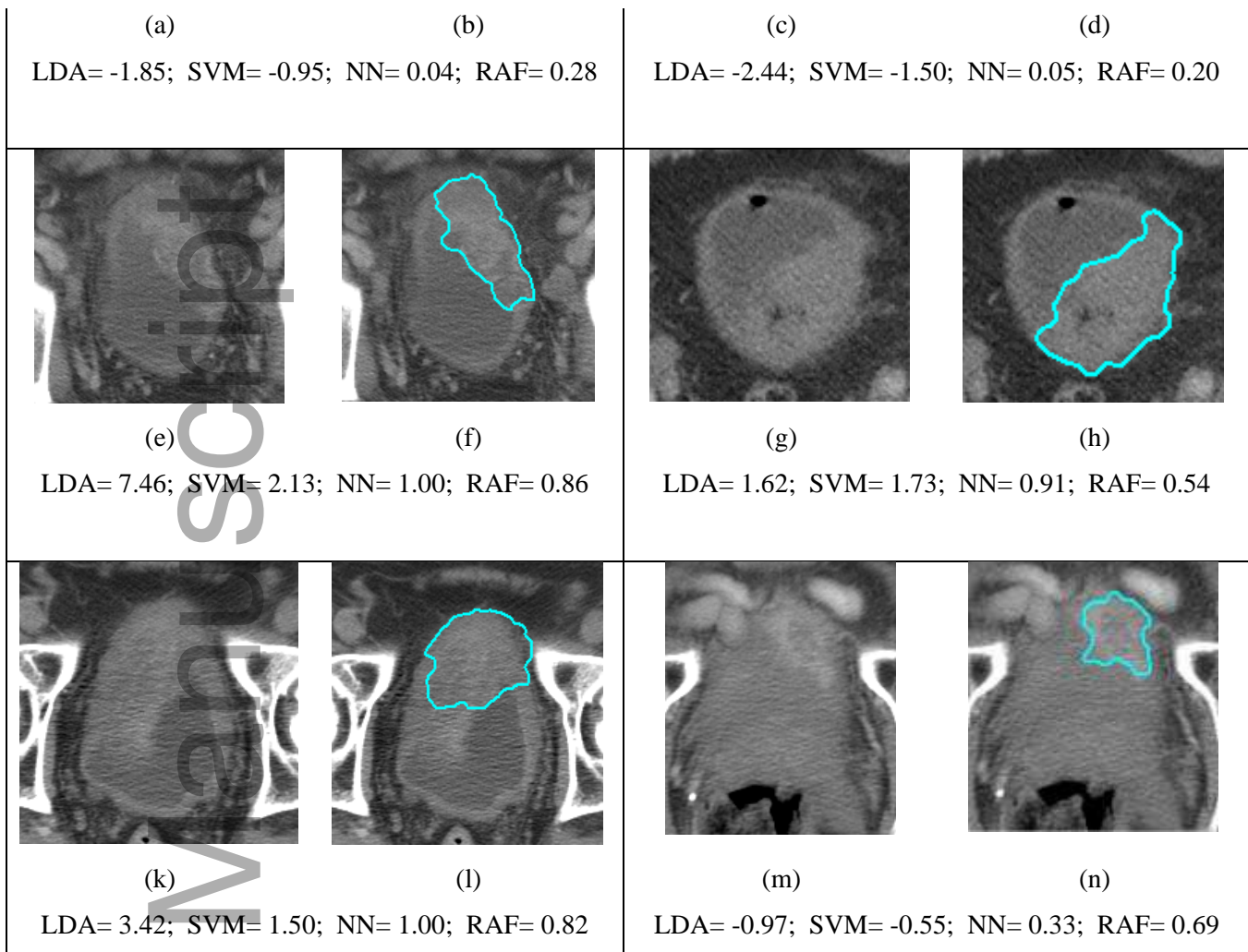


Figure 8. Examples of bladder cancers with stages \geq T2 or $<$ T2. The blue outlines represent the AI -CALS segmentation. The reported scores are test scores for the LDA, SVM, NN, and RAF classifiers based on the morphological features. Note that the output score ranges are different for different classifiers so that the score values should not be compared across classifiers. The two cases in (a)(b) and (c)(d) both contained was a T1 stage cancer that was properly classified with low scores from all classifiers. (e)(f) was a T3 stage case that was properly classified with high scores from all classifiers. (g)(h) was a T2 stage case that was properly classified with high scores from all classifiers. (k)(l) was a case that was clinically identified as T1 pre-surgery but was identified as a T2 stage cancer post-surgery. The classifiers classified the cancer as \geq T2 with high scores. (m)(n) was T2 stage cancer that was incorrectly identified by the LDA, SVM, and NN classifiers with low scores and correctly identified by the RAF with a high score.

377

378 We also have extracted features from the manually segmented bladder lesions and

379 applied the 4 different types of classifiers with the different feature sets to the cancer stage

380 prediction. The classifiers using features extracted from the manually segmented lesions
381 performed similarly to the classifiers using features extracted from the AI-CALS segmented
382 lesions. The test A_z values ranged from 0.77 to 0.95. For 6 out of the 24 experiments the
383 classifiers using features extracted from the manually segmented lesions performed better than
384 classifiers using features extracted from the AI-CALS segmentations. However, the differences
385 did not reach statistical significance. Therefore, although the performance of the AI-CALS lesion
386 segmentation was slightly lower than the radiologists' hand outlines the final classification
387 results were similar.

388 The main limitation of the study is the small data set. Another limitation is that we have
389 not applied the deep learning convolution neural network (DLCNN) to this bladder cancer
390 staging task. DLCNN has been shown to be superior to conventional classifiers in many
391 classification tasks, especially the classification of natural scene images with millions of training
392 samples. It also shows promise in number of medical imaging applications^{33,34} including bladder
393 segmentation³⁵ and bladder cancer treatment response monitoring³⁶. However, our experience
394 with DLCNN also indicates that it is not always the best, perhaps limited by the relatively small
395 annotated training set in medical imaging, even with transfer learning. As the performances of
396 the four conventional classifiers used in this study were quite high, it would not be a fair
397 comparison for DLCNN if we do not have adequate training for the latter. We will continue to
398 collect additional cases and compare the conventional classifiers with DLCNN for bladder
399 cancer staging in a future study.

400

401 6. CONCLUSION

402 In this preliminary study we proposed machine learning methods for prediction of
403 bladder cancer stage. It was found that the morphological features and texture features were
404 useful for assessing the stage of bladder lesions. The LDA, SVM, and NN classifiers all led to
405 relatively consistent results. There was a trend that the SVM and NN classifier slightly
406 outperformed the LDA classifier. The best overall results for the two-fold cross validation were
407 obtained when a combined feature subspace was used with the NN classifier. Further studies are
408 under way to improve the staging of bladder cancer and test the classifier on a larger data set,
409 and to investigate the potential of improving the predictive model by combining imaging
410 biomarkers with non-imaging biomarkers.

411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441

Acknowledgments

This work is supported by National Institutes of Health grant number U01CA179106.

References

¹American Cancer Society. *Cancer Facts & Figures 2017*. , (American Cancer Society, Inc., Atlanta, 2017).

²Bladder Cancer Advocacy Network, www.bcan.org/facts 2017, "Bladder Cancer Facts" " (2017).

³S. S. Chang, S. A. Boorjian, R. Chou, P. E. Clark, S. Daneshmand, B. R. Konety, R. Pruthi, D. Z. Quale, C. R. Ritch, J. D. Seigne, et al., "Diagnosis and Treatment of Non-Muscle Invasive Bladder Cancer: AUA/SUO Guideline," *Journal of Urology* 196, 1021-1029 (2016).

⁴J. A. Witjes, E. Comperat, N. C. Cowan, M. De Santis, G. Gakis, N. James, T. Le Bret, A. Sherif, A. G. Van der Heijden, and M. J. Ribal, "Guidelines on Muscle-invasive and Metastatic Bladder Cancer," *European Association of Urology* (2016).

⁵M. Babjuk, A. Bohle, M. Burger, E. Comperat, E. Kaasinen, J. Palou, M. Roupret, B. W. G. Van Rhijn, S. Shariat, R. Sylvester, et al., "Guidelines on Non-muscle-invasive Bladder Cancer (Ta, T1 and CIS)," *European Association of Urology* (2016).

⁶*AJCC Cancer Staging Handbook*, 8th ed. (American Joint Committee on Cancer, Chicago, IL, 2016).

⁷H. W. Herr and S. M. Donat, "Quality control in transurethral resection of bladder tumours," *Bju International* 102, 1242-1246 (2008).

- 442 ⁸J. J. Meeks, J. Bellmunt, B. H. Bochner, N. W. Clarke, S. Daneshmand, M. D. Galsky, N. M.
443 Hahn, S. P. Lerner, M. Mason, T. Powles, et al., "A Systematic Review of Neoadjuvant and
444 Adjuvant Chemotherapy for Muscle-invasive Bladder Cancer," *European Urology* 62, 523-533
445 (2012).
- 446
- 447 ⁹S. L. Fagg, P. Dawsonedwards, M. A. Hughes, T. N. Latief, E. B. Rolfe, and J. W. L. Fielding,
448 "CIS-Diamminedichloroplatinum (DDP) as initial treatment of invasive bladder cancer," *British*
449 *Journal of Urology* 56, 296-300 (1984).
- 450
- 451 ¹⁰D. Raghavan, B. Pearson, G. Coorey, W. Woods, D. Arnold, J. Smith, J. Donovan, and P.
452 Langdon, "Intravenous CIS-platinum for invasive bladder cancer – safety and feasibility of a
453 new approach," *Medical Journal of Australia* 140, 276-278 (1984).
- 454
- 455 ¹¹J. Huguet, M. Crego, S. Sabate, J. Salvador, J. Palou, and H. Villavicencio, "Cystectomy in
456 patients with high risk superficial bladder tumors who fail intravesical BCG therapy: Pre-
457 cystectomy prostate involvement as a prognostic factor," *European Urology* 48, 53-59 (2005).
- 458
- 459 ¹²H. M. Fritsche, M. Burger, R. S. Svatek, C. Jeldres, P. I. Karakiewicz, G. Novara, E. Skinner,
460 S. Denzinger, Y. Fradet, H. Isbarn, et al., "Characteristics and Outcomes of Patients with Clinical
461 T1 Grade 3 Urothelial Carcinoma Treated with Radical Cystectomy: Results from an
462 International Cohort," *European Urology* 57, 300-309 (2010).
- 463
- 464 ¹³P. Turker, P. J. Bostrom, M. L. Wroclawski, B. van Rhijn, H. Kortekangas, C. Kuk, T. Mirtti,
465 N. E. Fleshner, M. A. Jewett, A. Finelli, et al., "Upstaging of urothelial cancer at the time of
466 radical cystectomy: factors associated with upstaging and its effect on outcome," *Bju*
467 *International* 110, 804-811 (2012).
- 468
- 469 ¹⁴S. F. Shariat, G. S. Palapattu, P. I. Karakiewicz, C. G. Rogers, A. Vazina, P. J. Bastian, M. P.
470 Schoenberg, S. P. Lerner, A. I. Sagalowsky, and Y. Lotan, "Discrepancy between clinical and
471 pathologic stage: Impact on prognosis after radical cystectomy," *European Urology* 51, 137-151
472 (2007).

- 473
- 474 ¹⁵ACR *Manual on Contrast Media*, (ACR Committee on Drugs and Contrast Media, 2016).
- 475
- 476 ¹⁶L. M. Hadjiiski, H.-P. Chan, E. M. Caoili, R. H. Cohan, J. Wei, and C. Zhou, "Auto-Initialized
- 477 Cascaded Level Set (AI-CALS) Segmentation of Bladder Lesions on Multi-Detector Row CT
- 478 Urography," *Academic Radiology* 20, 148-155 (2013).
- 479
- 480 ¹⁷L. M. Hadjiiski, B. Sahiner, H.-P. Chan, N. Petrick, M. A. Helvie, and M. N. Gurcan, "Analysis
- 481 of Temporal Change of Mammographic Features: Computer-Aided Classification of Malignant
- 482 and Benign Breast Masses," *Medical Physics* 28, 2309-2317 (2001).
- 483
- 484 ¹⁸B. Sahiner, H.-P. Chan, N. Petrick, M. A. Helvie, and M. M. Goodsitt, "Computerized
- 485 characterization of masses on mammograms: The rubber band straightening transform and
- 486 texture analysis," *Medical Physics* 25, 516-526 (1998).
- 487
- 488 ¹⁹B. R. Dasarathy and E. B. Holder, "Image characterizations based on joint gray-level run-length
- 489 distributions," *Pattern Recog. Letters* 12, 497-502 (1991).
- 490
- 491 ²⁰T. W. Way, L. M. Hadjiiski, B. Sahiner, H.-P. Chan, P. N. Cascade, E. A. Kazerooni, N. Bogot,
- 492 and C. Zhou, "Computer-aided diagnosis of pulmonary nodules on CT scans: segmentation and
- 493 classification using 3D active contours," *Medical Physics* 33, 2323-2337 (2006).
- 494
- 495 ²¹H.-P. Chan, D. Wei, M. A. Helvie, B. Sahiner, D. D. Adler, M. M. Goodsitt, and N. Petrick,
- 496 "Computer-aided classification of mammographic masses and normal tissue: Linear discriminant
- 497 analysis in texture feature space," *Physics in Medicine and Biology* 40, 857-876 (1995).
- 498
- 499 ²²P. A. Lachenbruch, *Discriminant Analysis*, (Hafner Press, New York, 1975).
- 500
- 501 ²³M. M. Tatsuoka, *Multivariate Analysis, Techniques for Educational and Psychological*
- 502 *Research*, 2nd ed. (Macmillan, New York, 1988).
- 503

- 504 ²⁴D. E. Rumelhart, G. E. Hinton, and R. J. Williams, *Learning Internal Representation by Error*
505 *Propagation*, Parallel Distributed Processing (MIT Press, Cambridge, MA, 1986).
506
- 507 ²⁵V. N. Vapnik, *Statistical Learning Theory*, (Wiley, New York, 1998).
508
- 509 ²⁶C. J. C. Burges, "A Tutorial on Support Vector Machines for Pattern Recognition," *Data*
510 *Mining and Knowledge Discovery* 2, 121-167 (1998).
511
- 512 ²⁷T. K. Ho, "The random subspace method for constructing decision forests," *Ieee Transactions*
513 *on Pattern Analysis and Machine Intelligence* 20, 832-844 (1998).
514
- 515 ²⁸I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, *The WEKA Workbench. Online Appendix for*
516 *"Data Mining: Practical machine learning tools and techniques"*, (Morgan Kaufmann, 2016).
517
- 518 ²⁹P. Jaccard, "The distribution of the flora in the alpine zone," *New phytologist* 11, 37-50 (1912).
519
- 520 ³⁰C. E. Metz, "ROC methodology in radiologic imaging," *Investigative Radiology* 21, 720-733
521 (1986).
522
- 523 ³¹C. E. Metz, B. A. Herman, and J. H. Shen, "Maximum-likelihood estimation of receiver
524 operating characteristic (ROC) curves from continuously-distributed data," *Statistics in Medicine*
525 17, 1033-1053 (1998).
526
- 527 ³²"Metz ROC Software. University of Chicago Medical Center Department of Radiology,
528 see <http://metz-roc.uchicago.edu/MetzROC/software>,"
529
- 530 ³³G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghahfoorian, J. A. van der
531 Laak, B. van Ginneken, and C. I. Sánchez, "A Survey on Deep Learning in Medical Image
532 Analysis," arXiv:1702.05747 (2017).
533

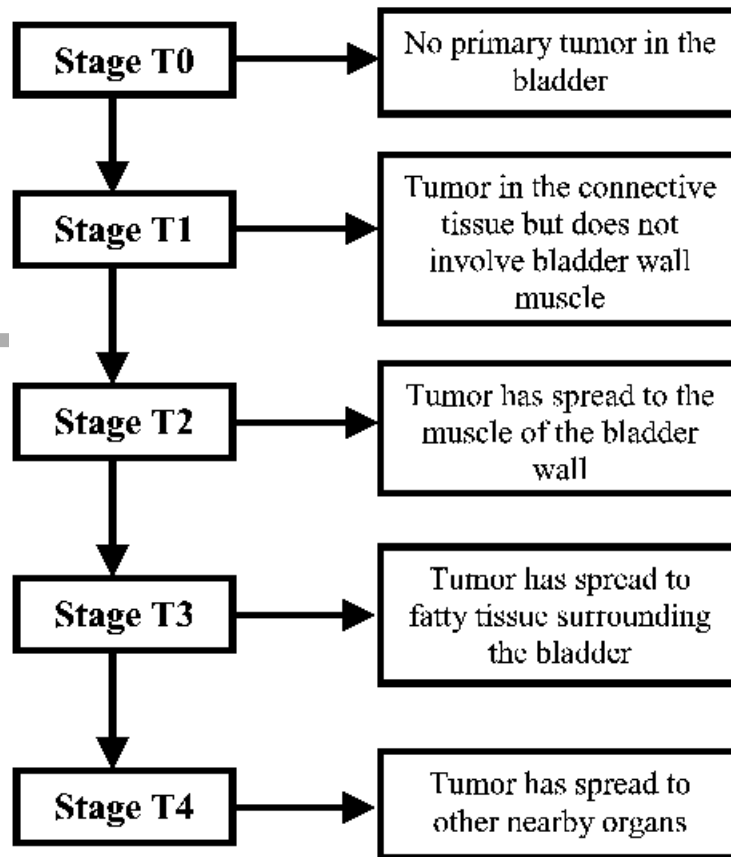
534 ³⁴H. Greenspan, B. van Ginneken, and R. M. Summers, "Deep Learning in Medical Imaging:
535 Overview and Future Promise of an Exciting New Technique," *Ieee Transactions on Medical*
536 *Imaging* 35, 1153-1159 (2016).

537
538 ³⁵K. H. Cha, L. Hadjiiski, R. K. Samala, H. P. Chan, E. M. Caoili, and R. H. Cohan, "Urinary
539 bladder segmentation in CT urography using deep-learning convolutional neural network and
540 level sets," *Medical Physics* 43, 1882-1896 (2016).

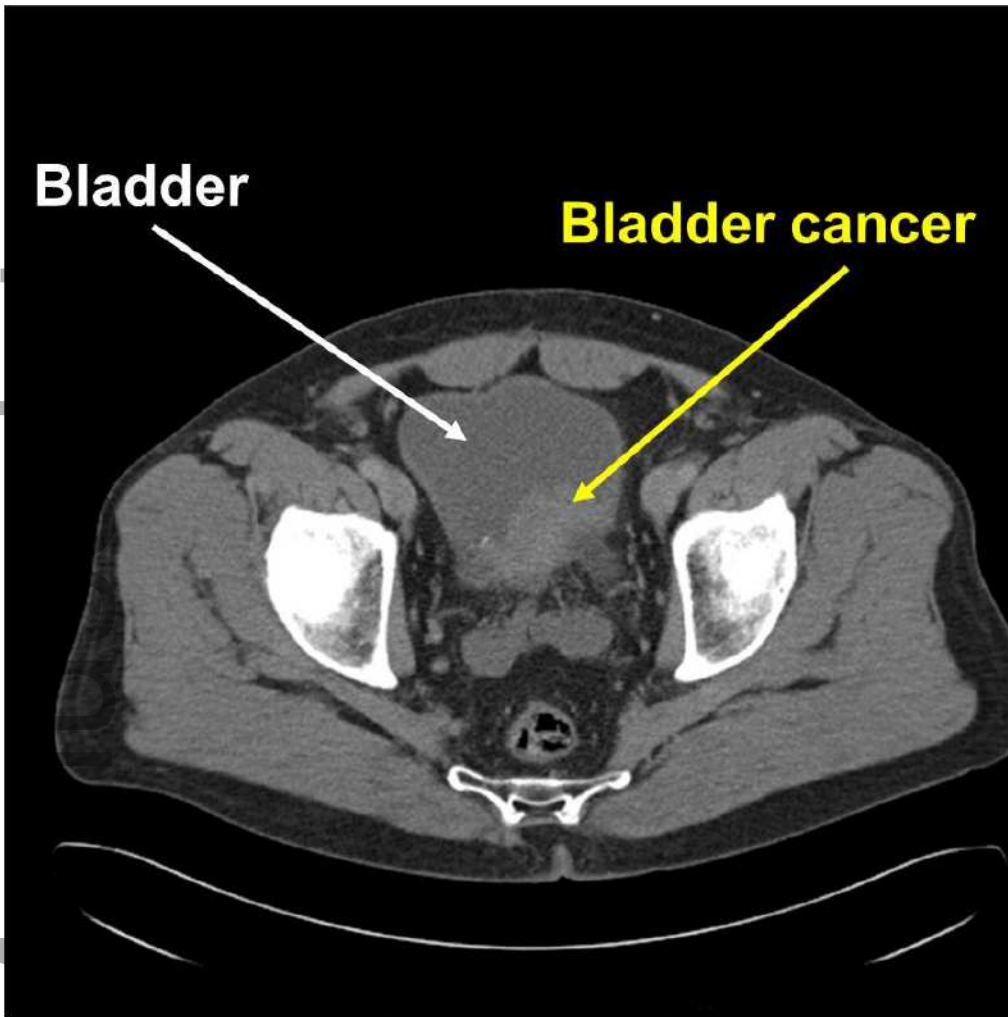
541
542 ³⁶K. H. Cha, L. M. Hadjiiski, H.-P. Chan, R. K. Samala, R. H. Cohan, E. M. Caoili, C.
543 Paramagul, A. Alva, and A. Z. Weizer, "Bladder cancer treatment response assessment using
544 deep learning in CT with transfer learning," *Proc SPIE* 10134, 101341-6 (2017).

545

Author Manuscript

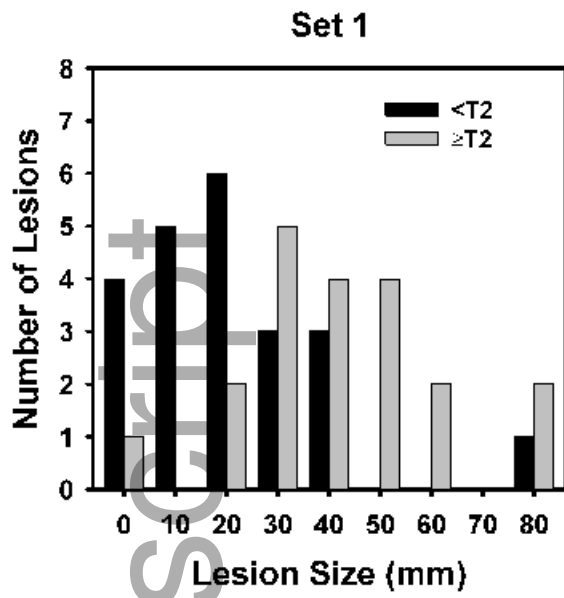


mp_12510_f1.tif

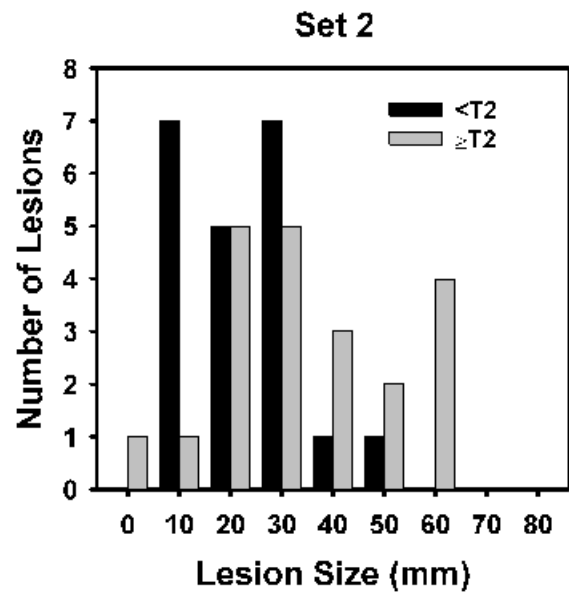


mp_12510_f2.tif

Author

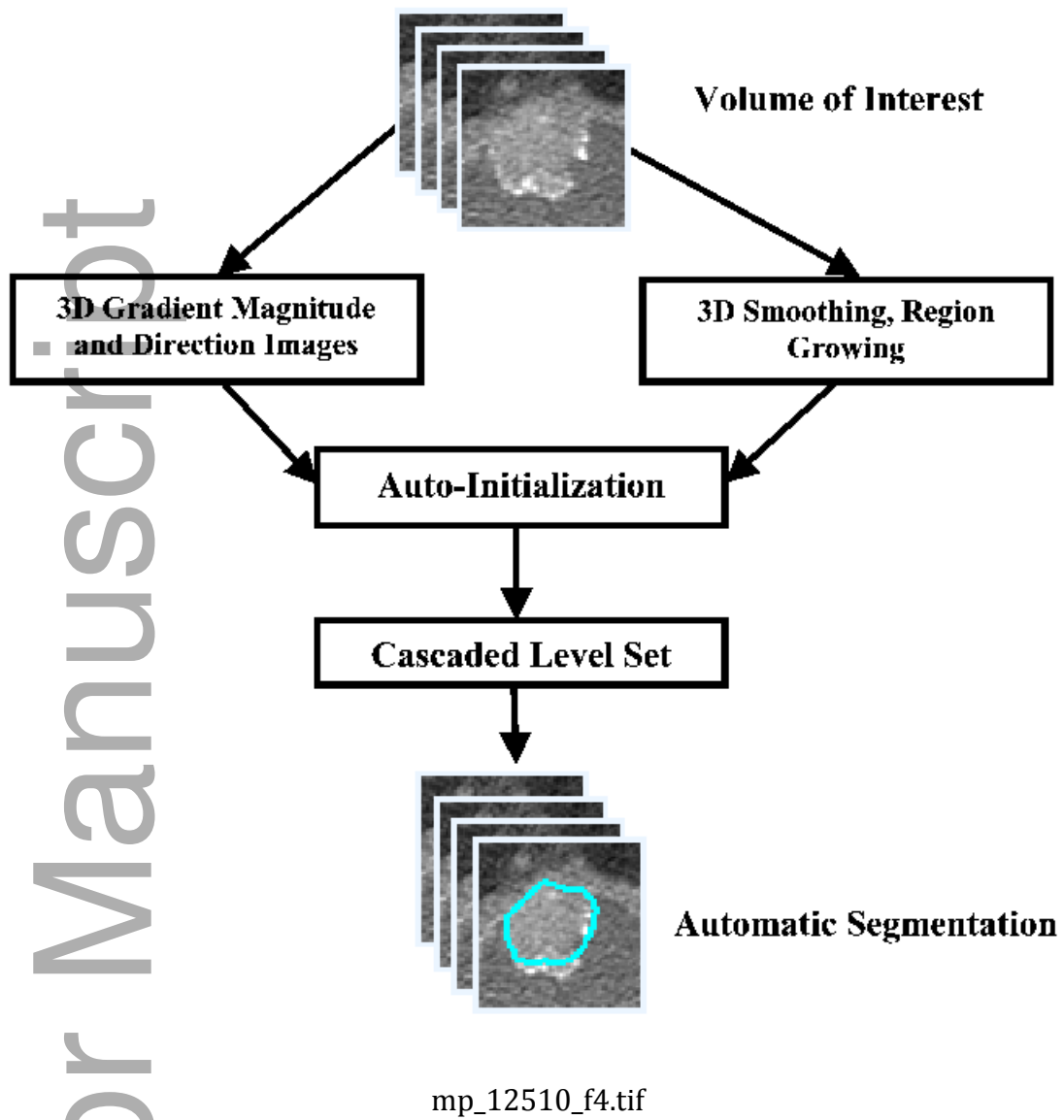


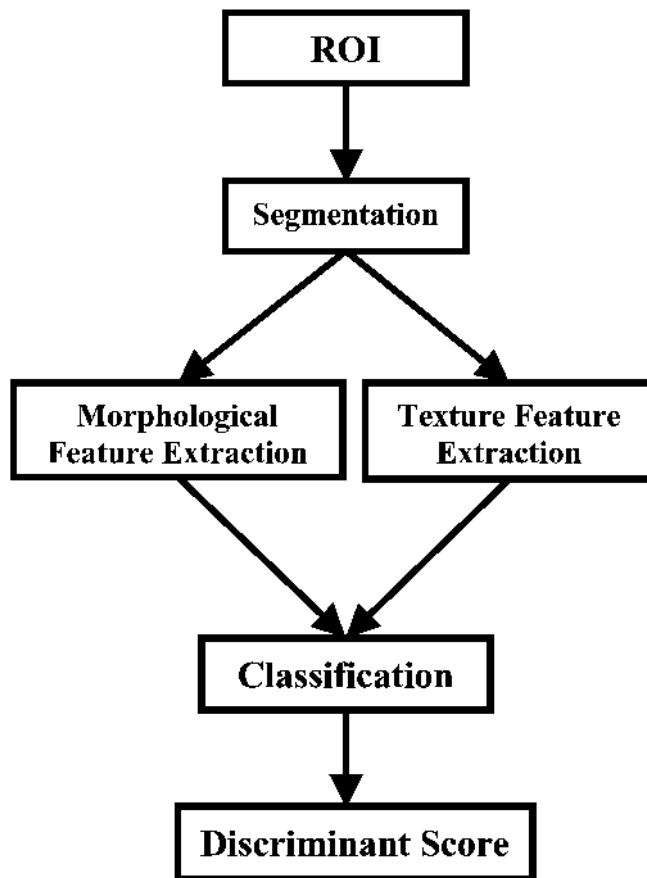
(a)



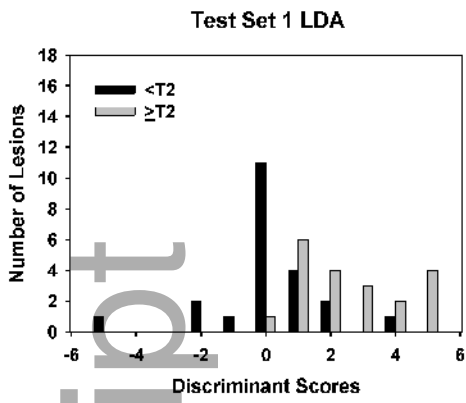
(b)

mp_12510_f3.tif

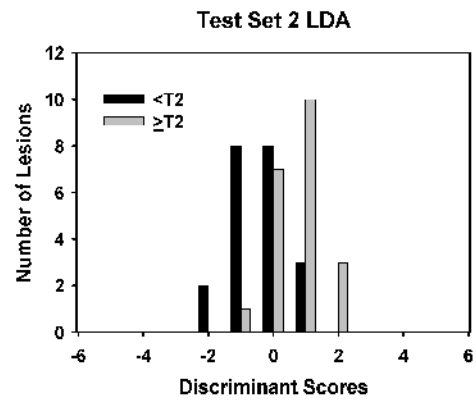




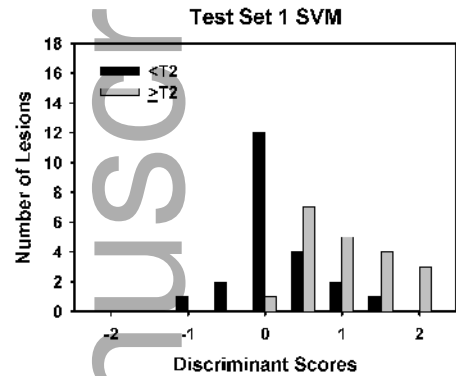
mp_12510_f5.tif



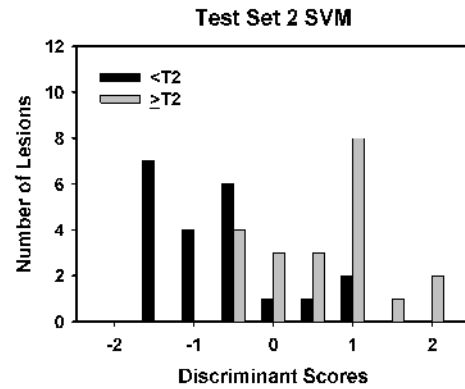
(a)



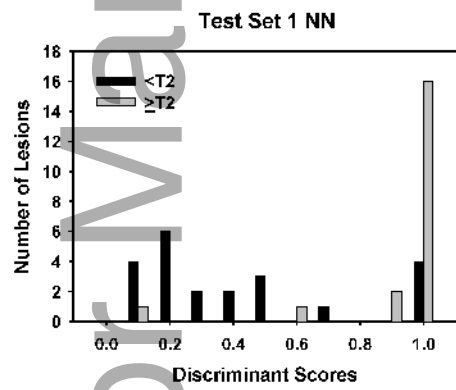
(b)



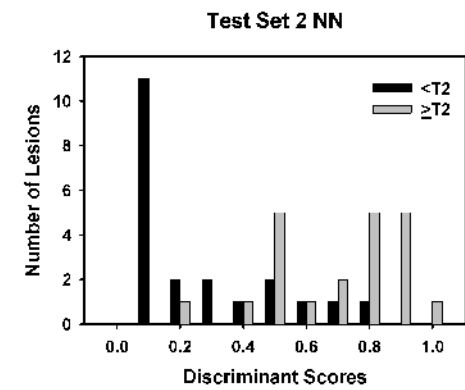
(c)



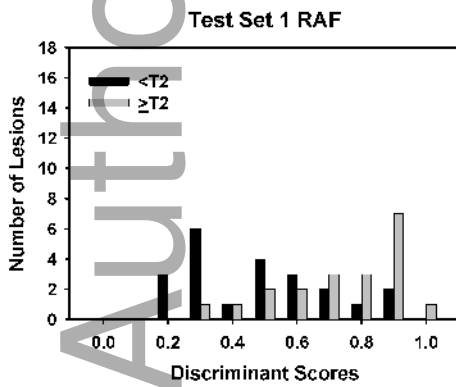
(d)



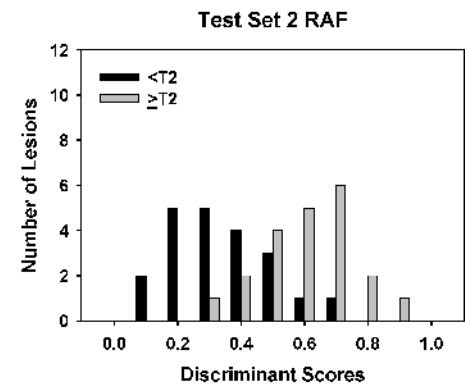
(e)



(f)

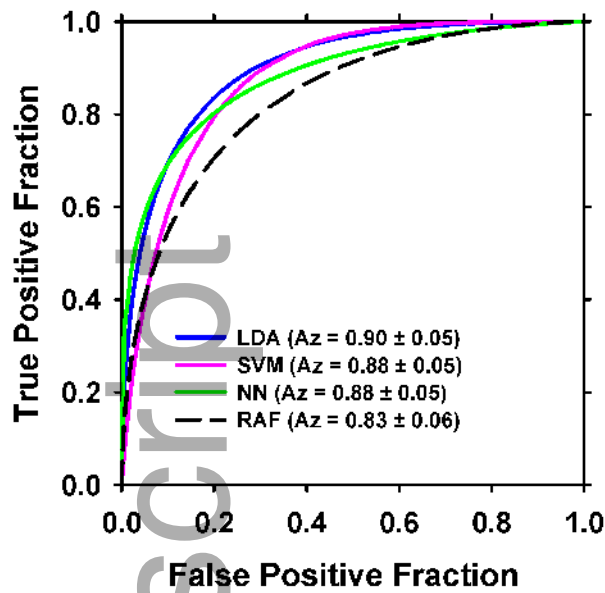


(g)

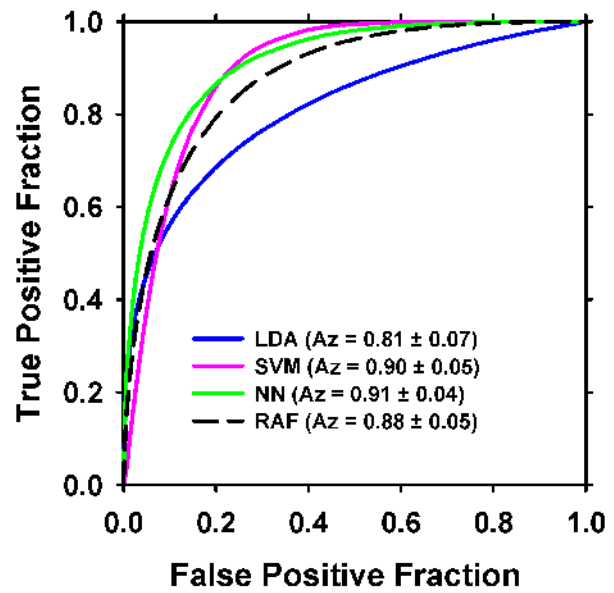


(h)

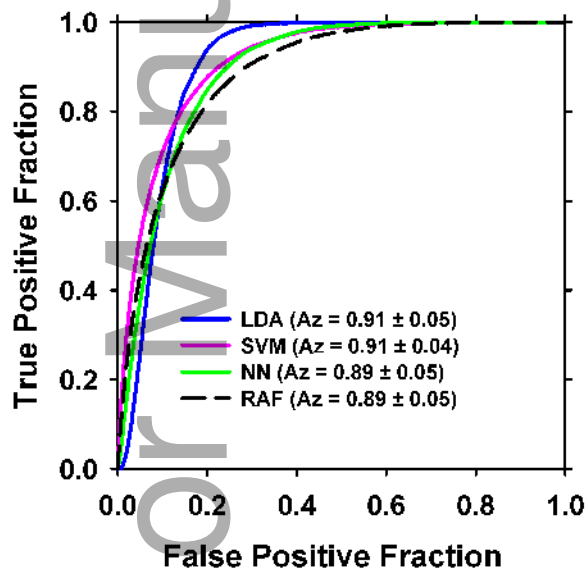
mp_12510_f6.tif



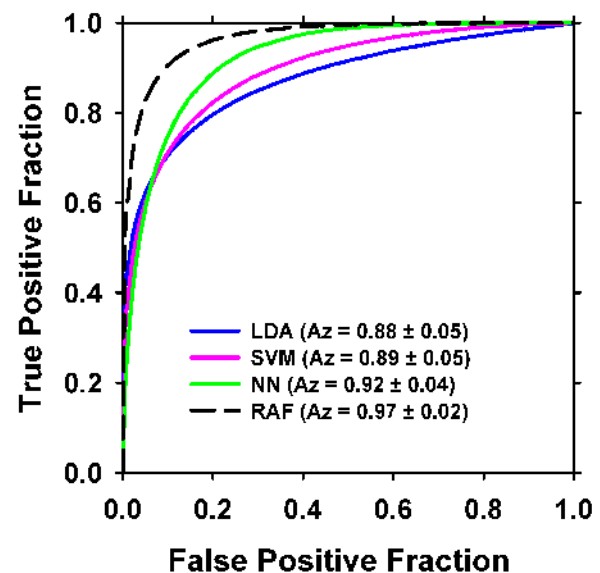
(a)



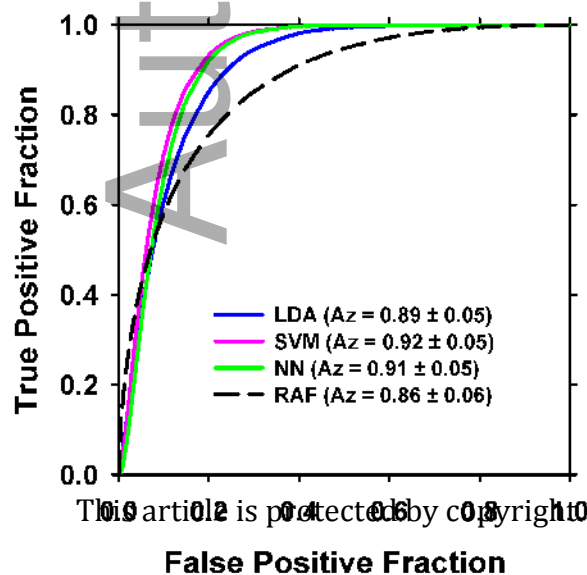
(b)



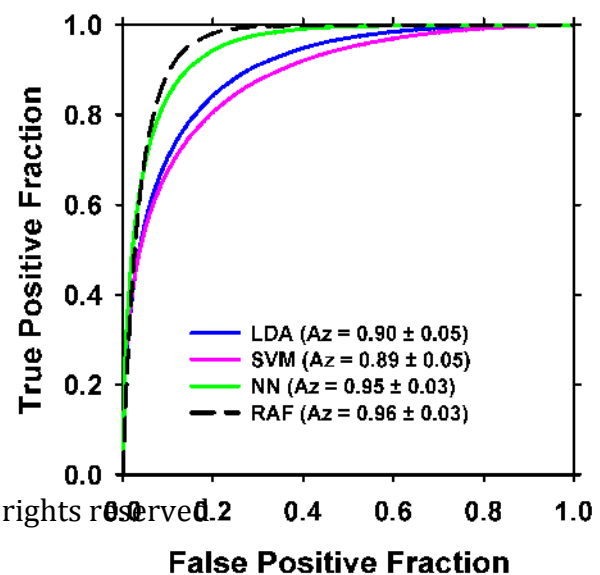
(c)



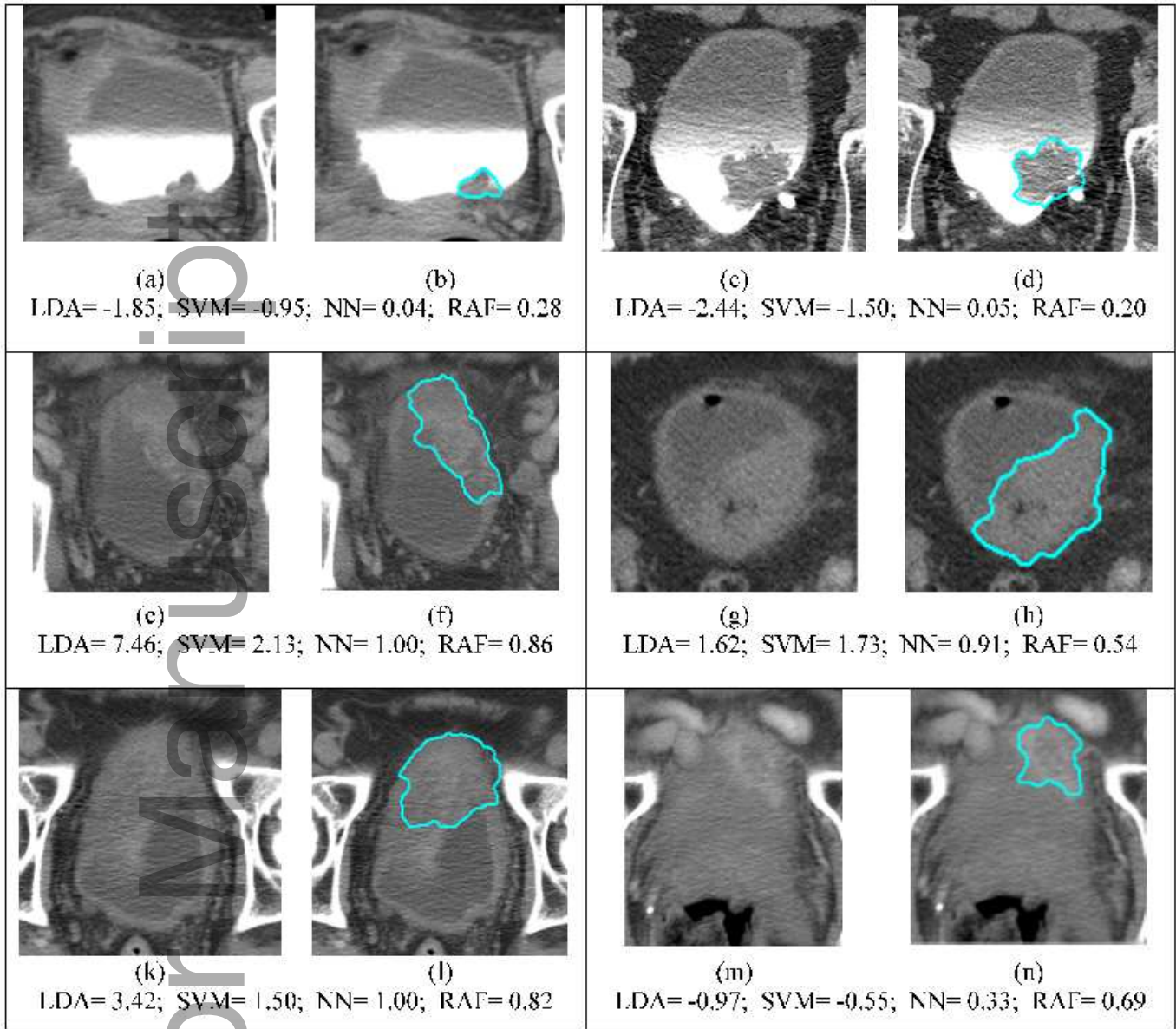
(d)



(e)



(f)



mp_12510_f8.tif