

Deep Learning Enables Prostate MRI Segmentation: A Large Cohort Evaluation with Inter-rater Variability Analysis

Yongkai Liu¹, Miao Qi^{1, 2*}, Chuthaporn Surawech^{1, 3}, Haoxin Zheng¹, Dan Nguyen⁴, Guang Yang⁵, Steven Raman¹, Kyung Hyun Sung¹

¹Department of Radiological Sciences, University of California, Los Angeles, United States, ²Department of Radiology, The First Affiliated Hospital of China Medical University, China, ³Department of Radiology, King Chulalongkorn Memorial Hospital, Thailand, ⁴Department of Radiation Oncology, University of Texas Southwestern Medical Center, United States, ⁵National Heart and Lung Institute, Imperial College London, United Kingdom

Submitted to Journal:
Frontiers in Oncology

Specialty Section:
Cancer Imaging and Image-directed Interventions

Article type:
Original Research Article

Manuscript ID:
801876

Received on:
25 Oct 2021

Revised on:
11 Nov 2021

Journal website link:
www.frontiersin.org

Conflict of interest statement

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest

Author contribution statement

Conceptualization, Y.L. and K.S.; methodology, Y.L.; software, Y.L.; validation, M.Q. and C.S.; formal analysis, Y.L. and K.S.; investigation, Y.L. and K.S.; resources, K.S. and S.R.; data curation, Y.L.; writing— original draft preparation, Y.L. and K.S.; writing—review and editing, Y.L., K.S., S.R., G.Y., Q.M., C.S. and D.N.; visualization, Y.L.; supervision, K.S. and S.R.; project administration, K.S.; funding acquisition, K.S. and G.Y.

Keywords

Prostate segmentation, deep attentive neural network, large cohort evaluation, qualitative evaluation, Quantitative evaluation, Volume measurement

Abstract

Word count: 264

Whole-prostate gland (WPG) segmentation plays a significant role in prostate volume measurement, treatment, and biopsy planning. This study evaluated a previously developed automatic WPG segmentation, deep attentive neural network (DANN), on a large, continuous patient cohort to test its feasibility in a clinical setting. With IRB approval and HIPAA compliance, the study cohort included 3,698 3T MRI scans acquired between 2016 and 2020. In total, 335 MRI scans were used to train the model, and 3,210 and 100 were used to conduct the qualitative and quantitative evaluation of the model. In addition, the DANN-enabled prostate volume estimation was evaluated by using 50 MRI scans in comparison with manual prostate volume estimation. For qualitative evaluation, visual grading was used to evaluate the performance of WPG segmentation by two abdominal radiologists, and DANN demonstrated either acceptable or excellent performance in over 96% of the testing cohort on the WPG or each prostate sub-portion (apex, midgland, or base). Two radiologists reached a substantial agreement on WPG and midgland segmentation ($\kappa=0.75$ and 0.63) and moderate agreement on apex and base segmentation ($\kappa=0.56$ and 0.60). For quantitative evaluation, DANN demonstrated a dice similarity coefficient of 0.93 ± 0.02 , significantly higher than other baseline methods, such as Deeplab v3+ and UNet (both p values < 0.05). For the volume measurement, 96% of the evaluation cohort achieved differences between the DANN-enabled and manual volume measurement within 95% limits of agreement. In conclusion, the study showed that the DANN achieved sufficient and consistent WPG segmentation on a large, continuous study cohort, demonstrating its great potential to serve as a tool to measure prostate volume.

Contribution to the field

The evaluation of current state-of-art deep learning methods was limited by relatively small sample size, ranging from tens to hundreds of MRI scans. It is relatively expensive to create large, continuous samples with manual segmentation of WPG, which limits the ability to test the DL models in a clinical setting. In this study, we evaluated a previously developed attentive deep learning-based automatic segmentation model using a large, continuous cohort of prostate 3T MRI scans ($n=3360$) to test its feasibility in a clinical setting.

Ethics statements

Studies involving animal subjects

Generated Statement: No animal studies are presented in this manuscript.

Studies involving human subjects

Generated Statement: The studies involving human participants were reviewed and approved by Institutional Review Board of UCLA. The ethics committee waived the requirement of written informed consent for participation.

Inclusion of identifiable human data

Generated Statement: No potentially identifiable human images or data is presented in this study.

Data availability statement

Generated Statement: The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

In review

Deep Learning Enables Prostate MRI Segmentation: A Large Cohort Evaluation with Inter-rater Variability Analysis

1 Yongkai Liu^{1,2}, Miao Qi^{1,6*}, Chuthaporn Surawech^{1,7}, Haoxin Zheng^{1,3}, Dan Nguyen⁵, Guang
2 Yang⁴, Steven S. Raman¹ and Kyunghyun Sung^{1,2}

3 ¹Department of Radiological Sciences, David Geffen School of Medicine, University of California,
4 Los Angeles, CA, USA

5 ²Physics and Biology in Medicine IDP, David Geffen School of Medicine, University of California,
6 Los Angeles, CA, USA

7 ³Department of Computer Science, Henry Samueli School of Engineering and Applied Science,
8 University of California, Los Angeles, CA, USA

9 ⁴National Heart and Lung Institute, Imperial College London, South Kensington, London, UK

10 ⁵Medical Artificial Intelligence and Automation Laboratory, Department of Radiation Oncology,
11 University of Texas Southwestern Medical Center, Dallas, TX, USA

12 ⁶Department of Radiology, The First Affiliated Hospital of China Medical University, Shenyang City
13 110001, Liaoning Province, China

14 ⁷Department of Radiology, Division of Diagnostic Radiology, Faculty of Medicine, Chulalongkorn
15 University and King Chulalongkorn Memorial Hospital, Bangkok, Thailand

16

17 * Correspondence:

18 Miao Qi

19 meganmiao89@gmail.com

20

21 Abstract

22 Whole-prostate gland (WPG) segmentation plays a significant role in prostate volume measurement,
23 treatment, and biopsy planning. This study evaluated a previously developed automatic WPG
24 segmentation, deep attentive neural network (DANN), on a large, continuous patient cohort to test its
25 feasibility in a clinical setting. With IRB approval and HIPAA compliance, the study cohort included
26 3,698 3T MRI scans acquired between 2016 and 2020. In total, 335 MRI scans were used to train the
27 model, and 3,210 and 100 were used to conduct the qualitative and quantitative evaluation of the model.
28 In addition, the DANN-enabled prostate volume estimation was evaluated by using 50 MRI scans in
29 comparison with manual prostate volume estimation. For qualitative evaluation, visual grading was
30 used to evaluate the performance of WPG segmentation by two abdominal radiologists, and DANN
31 demonstrated either acceptable or excellent performance in over 96% of the testing cohort on the WPG
32 or each prostate sub-portion (apex, midgland, or base). Two radiologists reached a substantial
33 agreement on WPG and midgland segmentation ($\kappa=0.75$ and 0.63) and moderate agreement on apex
34 and base segmentation ($\kappa=0.56$ and 0.60). For quantitative evaluation, DANN demonstrated a dice
35 similarity coefficient of 0.93 ± 0.02 , significantly higher than other baseline methods, such as Deeplab

36 v3+ and UNet (both p values < 0.05). For the volume measurement, 96% of the evaluation cohort
37 achieved differences between the DANN-enabled and manual volume measurement within 95% limits
38 of agreement. In conclusion, the study showed that the DANN achieved sufficient and consistent WPG
39 segmentation on a large, continuous study cohort, demonstrating its great potential to serve as a tool to
40 measure prostate volume.

41 **Keywords: prostate segmentation, deep attentive neural network, large cohort evaluation,**
42 **volume measurement, qualitative evaluation, quantitative evaluation**

43

44 1 Introduction

45 Whole-prostate gland (WPG) segmentation plays an important role in prostate volume measurement,
46 biopsy, and surgical planning [1]. Magnetic resonance imaging (MRI)-targeted transrectal ultrasound
47 fusion (MRI-fusion) biopsy has shown increased detection of clinically significant PCa and reduced
48 identification of clinically insignificant PCa [2], where the WPG segmentation is critical to enable the
49 MRI-fusion biopsy [3]. Also, prostate volume measurement via WPG segmentation can be used to
50 quantify the progression of benign prostatic hyperplasia [1] and to assist surgical planning [4].

51 Manual segmentation of WPG is time-consuming and laborious and commonly suffers from
52 inter-rater variability [5], making it inadequate for large-scale applications [6]. Deep learning (DL) [7–
53 10] has increasingly been utilized for the automatic segmentation of WPG. Zhu et al. [11] proposed a
54 deeply supervised convolutional neural network (CNN) using the convolutional information to
55 segment the prostate from MR images. Cheng et al. [8] developed a DL model with holistically nested
56 networks for prostate segmentation on MRI. Jia et al. [12] proposed an atlas registration and ensemble
57 deep CNN-based prostate segmentation. In addition, attentive DL [13] models were introduced to
58 enhance DL by paying attention to the particular regions of interest in an adaptive way and thus, have
59 achieved better segmentation performance than other DL-based models. However, to the best of our
60 knowledge, the evaluation of these methods was currently limited by relatively small sample size,
61 ranging from tens to hundreds of MRI scans. It is relatively expensive to create large, continuous
62 samples with manual segmentation of WPG, which limits the ability to test the DL models in a clinical
63 setting.

64 In this paper, we evaluated a previously developed DL-based automatic segmentation model,
65 deep attentive neural network (DANN) [13], using a large, continuous cohort of prostate 3T MRI scans
66 acquired between 2016 and 2020. The WPG segmentation by DANN was evaluated both quantitatively
67 and qualitatively. The quantitative evaluation was performed by using independent testing set with
68 manual segmentation as a ground-truth on a small dataset (n=100). The dice similarity coefficient
69 (DSC) [14] was used to measure the segmentation performance, compared with other baseline DL
70 methods. For qualitative evaluation, the segmentation performance was evaluated by two abdominal
71 radiologists independently via visual grading since the ground-truth manual segmentation was not
72 available for the large cohort (n=3,210). Inter-rater agreement between the two radiologists was
73 evaluated to check the consistency of the visual grading. We further investigated the segmentation on
74 different anatomical locations (i.e., apex, midgland, and base) as a secondary analysis. Finally, we
75 conducted the volume measurement using DANN-based segmentation on a small cohort (n=50)
76 (DANN-enabled volume measurement) and compared it with the manual volume measurement.

77 2 Materials and methods

78 2.1 MRI Datasets

79 With approval from the institutional review board (IRB), this retrospective study was carried out in
80 compliance with the United States Health Insurance Portability and Accountability Act (HIPAA) of
81 1996. After excluding MRI scans with severe artifacts and patients with prior surgery history and Foley
82 catheter, a total of 3,695 MRI scans, acquired on 3 T scanners (Skyra, Prisma, and Vida, Siemens
83 Healthineers, Erlangen, Germany), from January of 2016 to August of 2020, were included in the study.
84 Axial and coronal T2-weighted (T2W) Turbo spin-echo (TSE) images were used. Table 1 shows the
85 characteristics of the T2W MRI scan in the study.

86 Out of 3,695 3T MRI scans, 335 MRI scans (9%) were used as a training set, and the remaining
87 3,360 (91%) MRI scans were used as a testing set. **Training and testing datasets were randomly chosen**
88 **from the whole dataset.** The testing set included a qualitative evaluation set ($n=3,210$), a quantitative
89 evaluation set ($n=100$), and a volume measurement evaluation set ($n=50$). Table 2 shows the data
90 characteristics for each dataset. Training, quantitative, and volume measurement evaluation sets
91 required manual prostate contours as the segmentation reference standard. The manual annotation was
92 prepared by an abdominal radiologist (Q.M.) with more than five years of experience in the
93 interpretation of prostate MRI. In the training set, prostate contours were drawn on all axial T2W
94 images from all MRI scans, and on four mid-coronal T2W images (8th to 11th out of twenty slices) from
95 a subset of 100 MRI scans. In the quantitative and volume measurement evaluation sets, prostate
96 contours were drawn on all axial T2W images.

97 2.2 DL-based Whole Prostate Gland Segmentation Model

98 Figure 1 shows the overall workflow of the automatic WPG segmentation with DANN [13]. We added
99 the segmentation on the coronal plane to assist the selection of axial slices, reducing the inference time
100 of segmentation on the axial plane. During the testing, the workflow went through the following steps.
101 First, a $DANN_{cor}$, responsible for segmenting coronal slices, was adopted to segment the prostate on
102 the two-middle coronal images (9th and 10th slices out of twenty slices) for each MRI scan in the entire
103 testing set. The segmented coronal images were used to automatically select the axial T2W images that
104 contained the prostate gland. This would provide proper through-plane coverage of the prostate in the
105 axial slices. Next, $DANN_{ax}$ was used to perform the WPG segmentation on the selected axial T2W
106 images for each MRI scan in all the testing sets.

107 Both $DANN_{ax}$ and $DANN_{cor}$ were trained independently using the training set ($n=335$). First, a
108 subset of the training data ($n=100$) was used for training of $DANN_{cor}$, and four-middle coronal slices
109 (8th to 11th slices out of twenty slices) were used to make use of as many samples as possible. Once
110 the initial training of $DANN_{cor}$ was finished, two middle coronal slices were used as input to $DANN_{cor}$
111 for the rest of the training data. The segmented coronal slices by $DANN_{cor}$ were used to select certain
112 axial slices, and $DANN_{ax}$ was trained using all the selected axial slices in the entire training set.
113 Training and inferencing were conducted on a desktop computer with a 64-Linux system with 4 Titan
114 Xp GPU of 32 GB GDDR5 RAM. All the networks were trained with stochastic gradient descent as
115 the optimizer, with binary cross-entropy as the loss function. Pytorch was used to implement all the
116 DL networks. **The models were initially trained using 80% of the training dataset, and the rest of the**
117 **training dataset was used as the validation dataset. After the optimal hypermeters were found, we re-**

118 trained the models using the whole training dataset. The learning rate was initially set to 2.5e-3. All
 119 the networks were trained for 100 epochs with batch size 12.

120 2.3 Evaluation of Segmentation Performance

121 *Qualitative evaluation of segmentation performance*

122 We adopted the visual grading, similar to [15], to qualitatively evaluate the WPG segmentation. Two
 123 abdominal radiologists (M.Q. and C.S; each has over five years of experience in prostate MRI
 124 interpretation) assigned a visual grade, ranging from 1 to 3, to evaluate the segmentation performance,
 125 focusing on the whole prostate and sub-portions of the prostate (e.g., apex, midgland, and base). 1, 2,
 126 and 3 indicate unacceptable, acceptable, and excellent segmentation performance, respectively. Table
 127 3 shows the details when assigning the visual grade. Typical examples associated with each visual
 128 grade are shown in Figure 2. The readers independently ranked the segmentation quality. In addition,
 129 inter-rater reliability was assessed. To further investigate the segmentation at sub-portions of the
 130 prostate, we performed the sub-analysis for MRI scans without excellent segmentation performance
 131 agreed by both radiologists. Also, the segmentation performance for MRI scans with and without
 132 endorectal coil (ERC) was compared.

133 *Quantitative evaluation of segmentation performance*

134 3D DSC [16] was used to quantitatively evaluate and compare the segmentation performance in the
 135 quantitative evaluation set (n=100). The manual segmentations (M) were prepared by the radiologist
 136 on all axial slices as ground truths. DSC measures the overlapping between M and method-based (N)
 137 segmentation of the WPG volume and can be formulated as:

$$138 \quad DSC = \frac{2|M \cap N|}{|M \cup N|}, \quad (1)$$

139 where \cap and \cup indicate the intersection and union, respectively.

140 *Evaluation of volume measurement*

141 We further evaluated the performance of DANN-enabled volume measurements. After the radiologist
 142 manually drew the WPG contour on all axial slices, Pyradiomics [17] was used to calculate the prostate
 143 volume in the volume measurement evaluation set (n=50). The prostate volume from the DANN-based
 144 segmentation was compared with the manual volume measurement. The Bland-Altman plot [18] was
 145 used to analyze the agreement between manual and DANN-enabled WPG volume measurements.

146 2.4 Statistical Analysis

147 Mean and standard deviation were used to describe the distribution of DSC. The quantitative
 148 segmentation performance difference between the DANN and the baselines was compared using a
 149 paired sample t-test [19]. P values < 0.05 were considered statistically significant. Inter-rater reliability
 150 between two radiologists was measured by using the κ statistic [20]. The relationship between the
 151 value of κ and inter-rater reliability is listed as below, $\kappa < 0$: pool agreement; $0 < \kappa < 0.2$: slight agreement;
 152 $0.21 < \kappa < 0.4$: fair agreement; $0.41 < \kappa < 0.6$: moderate agreement; $0.61 < \kappa < 0.8$: substantial agreement;
 153 $0.81 < \kappa < 1.0$: almost perfect agreement.

154 3 RESULT

155 3.1 Qualitative Evaluation of WPG Segmentation

156 Figure 3 shows the proportion of acceptable or excellent segmentation quality in all MRI scans on the
 157 qualitative evaluation set at the whole prostate, or each sub-portion (apex, midgland, or base) of the
 158 prostate. The DANN method exhibited an acceptable or excellent segmentation performance in more
 159 than 96% of the MRI scans on the whole prostate or each sub-portion of the prostate. The segmentation
 160 at the midgland portion had achieved the best segmentation performance, while performed the worst
 161 at the base portion.

162 *Qualitative evaluation and inter-rater variability for WPG segmentation*

163 For WPG segmentation, 97.9% (n=3,141) and 93.2% (n=2,992) of the MRI scans were graded as
 164 having acceptable or excellent segmentation performance. Table 4 includes the confusion matrix to
 165 show the inter-rater variability of the visual grading. Overall, two readers reached a substantial
 166 consensus on the visual grading in 95.8% of the patients ($\kappa = 0.74$). When readers differed on the
 167 grading, the discrepancy in grading was less than one. 94.6% of segmentation results were unanimously
 168 considered as acceptable or excellent. Moreover, 91.5% of the MRI scans (n=2,861) were graded as
 169 having excellent segmentation performance according to the two radiologists. Unacceptable
 170 segmentation performance occurred only in 1.2% of the MRI scans (n=39), agreed by the two
 171 radiologists.

172 *Sub-analysis of MRI scans without excellent WPG segmentation*

173 We conducted the sub-analysis related to each sub-portion of the prostate (apex, midgland, or base)
 174 when the WPG segmentation was not excellent. The MRI scans with excellent segmentation agreed by
 175 two readers were excluded (n=2,929), and the rest of the MRI scans were used for the analysis (n=281).
 176 Figure 4 shows the confusion matrices of each sub-portion of the prostate on the rest of the MRI scans.
 177 46.3% of the MRI scans achieved the acceptable (or better) segmentation quality at the base slices,
 178 while 94.3% and 83.3% of the MRI scans achieved the acceptable (or better) segmentation quality at
 179 the midgland and apex slices.

180 *Comparison between MRI scans with and without ERC*

181 We compared the WPG segmentation quality for the MRI scans with and without ERC [21]. Figure 5
 182 shows the confusion matrices of the visual grades of segmentation on MRI scans with and without
 183 ERC. There were substantial agreements ($\kappa = 0.64$ and 0.85) between the two radiologists on WPG
 184 segmentation of MRI scans with and without ERC. When considering the inter-rater agreement of
 185 WPG segmentation, DANN demonstrated acceptable WPG performance in more than 95.5% of MRI
 186 scans with ERC compared to 84.3% of those without ERC. MRI scans with ERC had a larger
 187 proportion of unacceptable WPG segmentation compared to those without ERC (12.1% vs. 2.2%).

188 3.2 Quantitative Evaluation of WPG Segmentation

189 The quantitative performance of the DANN was compared to the other two baseline methods, including
 190 Deeplab v3+ [22] and UNet [23]. Table 5 shows the comparisons of DSCs between DANN and the
 191 baseline methods. The DANN achieved a DSC of 0.93, which was higher than those of Deeplab v3+
 192 and UNet with significant differences (both p values < 0.05).

193 3.3 Evaluation of Volume Measurement

194 Figure 6 shows the agreement between manual and DANN-enabled volume measurements in the
195 Bland-Altman plot. The mean difference between the two-volume measurements was calculated as an
196 estimated bias. Standard deviation (SD) of the differences and 95% limits of agreement (average
197 difference ± 1.96 SD) were calculated to assess the random fluctuations around this mean. 48 out of
198 50 cases (96%) had the volume measurement differences within 95% limits of agreement, indicating
199 that the manual and DANN-enabled volume measurements can be potentially used interchangeably.

200 4 DISCUSSION

201 A deep attentive neural network [13], DANN, for the automatic WPG segmentation was evaluated on
202 a large, continuous patient cohort. In the qualitative evaluation, DANN demonstrated that the
203 segmentation quality is either acceptable or excellent in most cases. Two radiologists exhibited a
204 substantial agreement for the qualitative evaluation. In the quantitative evaluation, DANN exhibited a
205 significantly higher DSC than the baseline methods, such as UNet and Deeplab v3+. Also, 96% of the
206 testing cohort had volume measurement differences within 95% limits of agreement.

207 We found that DANN demonstrated worse segmentation performance at the prostate base than
208 at the apex and midgland slices. This may be due to the fact that the anatomical structure of the prostate
209 base is relatively more complex than other prostate portions. The prostate base is in continuity with the
210 bladder and seminal vesicles, and thus the boundary may contain partial volume effects and mild
211 movement artifacts.

212 We observed that the segmentation performance was somewhat limited when MRI scans were
213 acquired with an ERC. We believe that this may be because there were only three MRI scans with ERC
214 in the training dataset. A large training data with ERC may allow the model to learn representative
215 features related to the prostate MRI with ERC. In addition, images often exhibit large intensity
216 variation compared to the MRI scans without ERC as ERC is close to the prostate. This may require
217 including an even larger training dataset to account for these intensity variations than those without
218 ERC.

219 We refined DANN by adding the coronal segmentation to assist the selection of axial slices for
220 WPG segmentation. With assistance from the coronal segmentation, the axial model conducted the
221 segmentation only on the selected axial slices instead of applying it to all axial slices, which reduces
222 the inference time. Table 6 contains the inference time between the segmentation with and without
223 coronal segmentation. The total inference time in a combination of coronal and axial slices was 25%
224 less than the inference time without assisting the selection of axial slices (12.6 min vs. 16.4 min). In
225 addition, we observed that DSC was not different when adding the coronal segmentation in the
226 quantitative evaluation.

227 Compared with quantitative evaluation, qualitative evaluation includes unique characteristics
228 and benefits. The DSC-based evaluation often overlooks the segmentation performance on small
229 regions when they were combined with larger regions. Prostate at apex or base slices is relatively
230 smaller than the one in the middle, and therefore, the quantitative evaluation may not be sensitive
231 enough to illustrate limitations at these locations when 3D DSC is used for the evaluation. Also, the
232 DSC-based evaluation is not directly associated with clinical implications, while qualitative evaluation
233 categorized the segmentation results based on the quality to which segmentation can be acceptable
234 clinically.

235 Our study still has a few limitations: 1) the MRI scans in this study were acquired from three
236 3T MRI scanners at a single medical center. Prostate MRI sequence parameters are generally well-
237 standardized by the Prostate Imaging–Reporting and Data System (PI-RADS) guidelines [24], but
238 future studies would include testing DANN at multiple institutions. 2) the inter-rater variability was
239 tested between two radiologists. We will include more radiologists to evaluate comprehensive inter-
240 rater variability. 3) large GPU memory was required during the training and testing since DANN
241 included the spatial attention mechanism that caused considerable computational complexity. In the
242 future, we will explore the criss-cross attention module [25] that uses the contextual information of all
243 the pixels on the criss-cross path for each pixel, which has shown the potential to reduce the GPU
244 memory.

245

246 **5 Conclusion**

247 Our study showed that the proposed deep learning-based prostate segmentation (DANN) could
248 generate segmentation of the prostate with sufficient quality in a consistent manner when a large,
249 continuous cohort of prostate MRI scans was used for evaluation. The qualitative evaluation conducted
250 by two abdominal radiologists showed that 95% of the segmentation results were either acceptable or
251 excellent with a great inter-rater agreement. In the quantitative evaluation, DANN was superior to the
252 state-of-art deep learning methods, and the difference between manual and DANN-enabled volume
253 measurements was subtle in most cases. The method has a great potential to serve as a tool to assist
254 prostate volume measurements, and biopsy and surgical planning in a clinically relevant setting.

255 **6 Funding**

256 This work was supported in part by the National Institutes of Health R01-CA248506 and funds from
257 the Integrated Diagnostics Program, Departments of Radiological Sciences and Pathology, David
258 Geffen School of Medicine, UCLA; This study was also supported in part by the British Heart
259 Foundation (Project Number: TG/18/5/34111, PG/16/78/32402), the European Research Council
260 Innovative Medicines Initiative (DRAGON, H2020-JTI-IMI2 101005122), the AI for Health Imaging
261 Award (CHAIMELEON, H2020-SC1-FA-DTS-2019-1 952172), and the UK Research and Innovation
262 Future Leaders Fellowship (MR/V023799/1).

263

264

265 **7 References**

- 266 1. Garvey B, Türkbey B, Truong H, Bernardo M, Periaswamy S, Choyke PL (2014) Clinical value of prostate
 267 segmentation and volume determination on MRI in benign prostatic hyperplasia. *Diagnostic Interv*
 268 *Radiol* 20:229.
- 269 2. Ahmed HU, Bosaily AE-S, Brown LC, Gabe R, Kaplan R, Parmar MK, Collaco-Moraes Y, Ward K,
 270 Hindley RG, Freeman A, others (2017) Diagnostic accuracy of multi-parametric MRI and TRUS biopsy
 271 in prostate cancer (PROMIS): a paired validating confirmatory study. *Lancet* 389:815–822.
- 272 3. Kasivisvanathan V, Rannikko AS, Borghi M, Panebianco V, Mynderse LA, Vaarala MH, Briganti A,
 273 Budäus L, Hellawell G, Hindley RG, others (2018) MRI-targeted or standard biopsy for prostate-cancer
 274 diagnosis. *N Engl J Med* 378:1767–1777.
- 275 4. Oelke M, Bachmann A, Descazeaud A, Emberton M, Gravas S, Michel MC, N’dow J, Nordling J, Jean J
 276 (2013) EAU guidelines on the treatment and follow-up of non-neurogenic male lower urinary tract
 277 symptoms including benign prostatic obstruction. *Eur Urol* 64:118–140.
- 278 5. Wenger E, Mårtensson J, Noack H, Bodammer NC, Kühn S, Schaefer S, Heinze H-J, Düzel E, Bäckman L,
 279 Lindenberger U, others (2014) Comparing manual and automatic segmentation of hippocampal volumes:
 280 reliability and validity issues in younger and older brains. *Hum Brain Mapp* 35:4236–4248.
- 281 6. Yuan Y, Li B, Meng MQ-H (2015) Bleeding frame and region detection in the wireless capsule endoscopy
 282 video. *IEEE J Biomed Heal informatics* 20:624–630.
- 283 7. Jin Y, Yang G, Fang Y, Li R, Xu X, Liu Y, Lai X (2021) 3D PBV-Net: An automated prostate MRI data
 284 segmentation method. *Comput Biol Med* 128:104160.
- 285 8. Cheng R, Roth HR, Lay NS, Lu L, Turkbey B, Gandler W, McCreedy ES, Pohida TJ, Pinto PA, Choyke
 286 PL, others (2017) Automatic magnetic resonance prostate segmentation by deep learning with
 287 holistically nested networks. *J Med imaging* 4:41302.
- 288 9. Checcucci E, Autorino R, Cacciamani GE, Amparore D, De Cillis S, Piana A, Piazzolla P, Vezzetti E, Fiori
 289 C, Veneziano D, others (2019) Artificial intelligence and neural networks in urology: current clinical
 290 applications. *Minerva Urol e Nefrol Ital J Urol Nephrol* 72:49–57.
- 291 10. Checcucci E, De Cillis S, Granato S, Chang P, Afyouni AS, Okhunov Z, others (2020) Applications of
 292 neural networks in urology: a systematic review. *Curr Opin Urol* 30:788–807.
- 293 11. Zhu Q, Du B, Turkbey B, Choyke PL, Yan P (2017) Deeply-supervised CNN for prostate segmentation.
 294 2017 Int. Jt. Conf. neural networks. pp 178–184
- 295 12. Jia H, Xia Y, Song Y, Cai W, Fulham M, Feng DD (2018) Atlas registration and ensemble deep
 296 convolutional neural network-based prostate segmentation using magnetic resonance imaging.
 297 *Neurocomputing* 275:1358–1369.
- 298 13. Liu Y, Yang G, Hosseiny M, Azadikhah A, Mirak SA, Miao Q, Raman SS, Sung K (2020) Exploring
 299 Uncertainty Measures in Bayesian Deep Attentive Neural Networks for Prostate Zonal Segmentation.
 300 *IEEE Access* 8:151817–151828.
- 301 14. Dice LR (1945) Measures of the Amount of Ecologic Association Between Species. *Ecology* 26:297–302.

- 302 15. Kiri\esli HA, Schaap M, Klein S, Papadopoulou S-L, Bonardi M, Chen C-H, Weustink AC, Mollet NR,
303 Vonken E-J, van der Geest RJ, others (2010) Evaluation of a multi-atlas based method for segmentation
304 of cardiac CTA data: a large-scale, multicenter, and multivendor study. *Med Phys* 37:6279–6291.
- 305 16. Liu Y, Yang G, Mirak SA, Hosseiny M, Azadikhah A, Zhong X, Reiter RE, Lee Y, Raman SS, Sung K
306 (2019) Automatic Prostate Zonal Segmentation Using Fully Convolutional Network With Feature
307 Pyramid Attention. *IEEE Access* 7:163626–163632.
- 308 17. Kitzing YX, Prando A, Varol C, Karczmar GS, Maclean F, Oto A (2016) Benign conditions that mimic
309 prostate carcinoma: MR imaging features with histopathologic correlation. *Radiographics* 36:162–175.
- 310 18. Giavarina D (2015) Understanding bland altman analysis. *Biochem medica* 25:141–151.
- 311 19. Semenick D (1990) Tests and measurements: The T-test. *Strength \& Cond J* 12:36–37.
- 312 20. McHugh ML (2012) Interrater reliability: the kappa statistic. *Biochem medica* 22:276–282.
- 313 21. Turkbey B, Merino MJ, Gallardo EC, Shah V, Aras O, Bernardo M, Mena E, Daar D, Rastinehad AR,
314 Linehan WM, others (2014) Comparison of endorectal coil and nonendorectal coil T2W and diffusion-
315 weighted MRI at 3 Tesla for localizing prostate cancer: correlation with whole-mount histopathology. *J*
316 *Magn Reson Imaging* 39:1443–1448.
- 317 22. Chen LC, Zhu Y, Papandreou G, Schroff F, Adam H (2018) Encoder-decoder with atrous separable
318 convolution for semantic image segmentation. *Lect. Notes Comput. Sci. (including Subser. Lect. Notes*
319 *Artif. Intell. Lect. Notes Bioinformatics)*. pp 833–851
- 320 23. Ronneberger O, Fischer P, Brox T (2015) U-net: Convolutional networks for biomedical image
321 segmentation. *Int. Conf. Med. image Comput. Comput. Interv.* pp 234–241
- 322 24. Turkbey B, Rosenkrantz AB, Haider MA, Padhani AR, Villeirs G, Macura KJ, Tempany CM, Choyke PL,
323 Cornud F, Margolis DJ, others, Thoeny HC, Verma S, Barentsz J, Weinreb JC (2019) Prostate Imaging
324 Reporting and Data System Version 2.1: 2019 Update of Prostate Imaging Reporting and Data System
325 Version 2. *Eur Urol.* doi: 10.1016/j.eururo.2019.02.033
- 326 25. Huang Z, Wang X, Huang L, Huang C, Wei Y, Liu W (2019) Ccnet: Criss-cross attention for semantic
327 segmentation. *Proc. IEEE/CVF Int. Conf. Comput. Vis.* pp 603–612

328

329

330

331

332

333 **TABLES**

334 **Table 1:** T2-weighted TSE MRI sequence parameters in the study.

View	Axial	Coronal
Matrix size	320 × 320	320 × 320
Flip angle	160°	147°
Resolution	0.625 × 0.625 × 3.6	0.625 × 0.625 × 3.6
Field of View (mm ²)	200 × 200	200 × 200
Repetition Time (ms)	3000-7480	2880-7200
Echo Time (ms)	97-112	97-109
Number of slices	20	20
Scan Time (s)	200	200

ms: Millisecond; s: second; mm: millimeter;

335

336

337

338

339

340

341 **Table 2:** Data characteristics in the training, qualitative, and quantitative evaluation.

		Training Dataset	Qualitative Evaluation Dataset	Quantitative Evaluation Dataset	Volume Evaluation Dataset
Number of MRI scans		335	3,210	100	50
Number of patients with Endo-Rectal Coil		3	84	0	0
MRI scans with different vendors	Skyra	295	2,806	93	45
	Prisma	10	145	4	3
	Vida	30	259	3	2

342

343

344

345

346

347

348

349

350 **Table 3:** Description of each visual grade for qualitative segmentation evaluation.

Score	Visual scoring description
3	The segmentation is excellent. The vast majority (>90%) of the prostate region has been correctly segmented and the percentage of prostate slices with the failure segmentation is less than 10%.
2	The segmentation is acceptable. Most of the region (>70%) is correctly segmented, and the percentage of prostate slices that the method fails to segment is less than 30%.
1	The segmentation is unacceptable. More than 30% of the prostate region has been not correctly segmented or wrongly segmented, and the percentage of prostate slices that the method fails to segment is larger than 30%.

351

352

353

354

355

356

357

358

359

360 **Table 4:** Confusion matrices between the visual grades assigned by two readers. Kappa coefficient
 361 (κ) is used to measure the inter-rater variability between the two readers.

All	Reader 2			Kappa (κ)	
	Visual grade	1	2	3	
Reader 1	1	47 (1.5)	1 (0.0)	0 (0.0)	Substantial agreement ($\kappa = 0.75$)
	2	22 (0.7)	99 (3.1)	49 (1.5)	
	3	0 (0.0)	63 (2.0)	2,929 (91.3)	

362

363

364

365

366

367

368

369

370

371

372

In review

373

Table 5. Quantitative DSC comparisons with baseline methods

Methods	DSC
Proposed Method	0.93±0.02
Deeplab v3+	0.92±0.02 P<0.05
UNet	0.91±0.03 P<0.05

374

375

376

377

378

379

380

381

382

In review

383 **Table 6.** Inference time estimation and DSCs obtained with and without coronal segmentation
 384 assistance

	Without coronal segmentation assistance	With coronal segmentation assistance
Overall inference time estimation in the qualitative evaluation	16.4 minutes (67,775)	12.6 minutes (45,713)
DSCs obtained in the quantitative evaluation	0.93	0.93

() indicates the total amount of MRI slices the method needed to segment.

385

386

387

388

389

390

391

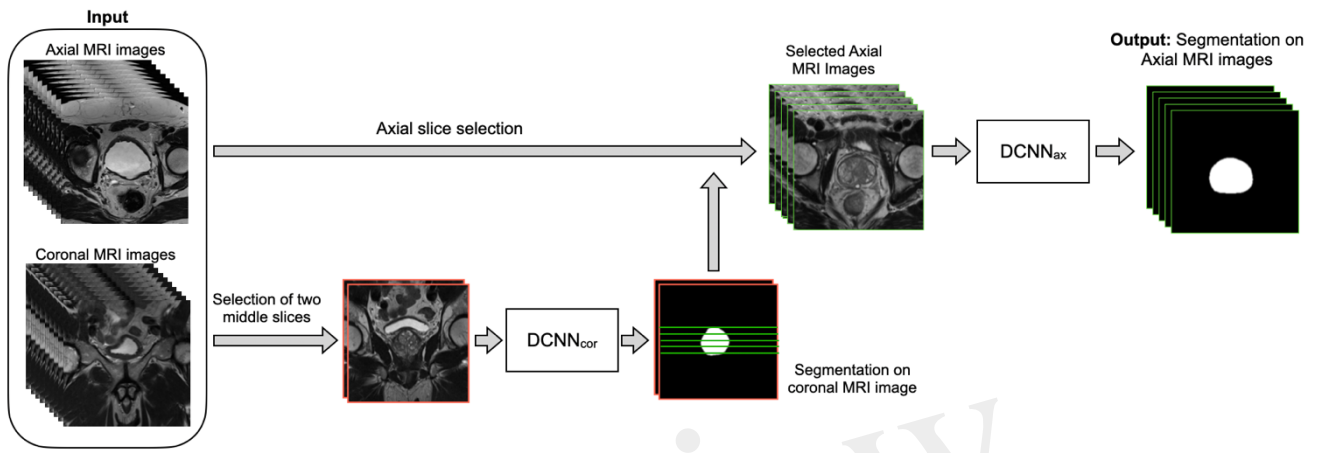
392

393

394

In review

395 **FIGURE CAPTIONS**



396

397 **Figure 1:** The overall workflow of the automatic WPG segmentation with DANN. Both axial and
 398 coronal T2W images were used as input, where the coronal images were used to assist the selection of
 399 certain axial images containing the prostate gland. DANN_{cor} was firstly performed on the two middle
 400 coronal images, indicated by images with the red border. Next, green lines selected by the prostate
 401 segmentation on the coronal images were used to determine the selection of axial slices (images with
 402 green borders). Once the axial images were selected, DANN_{ax} was performed on the axial MRI slices
 403 for the segmentation of WPG.

404

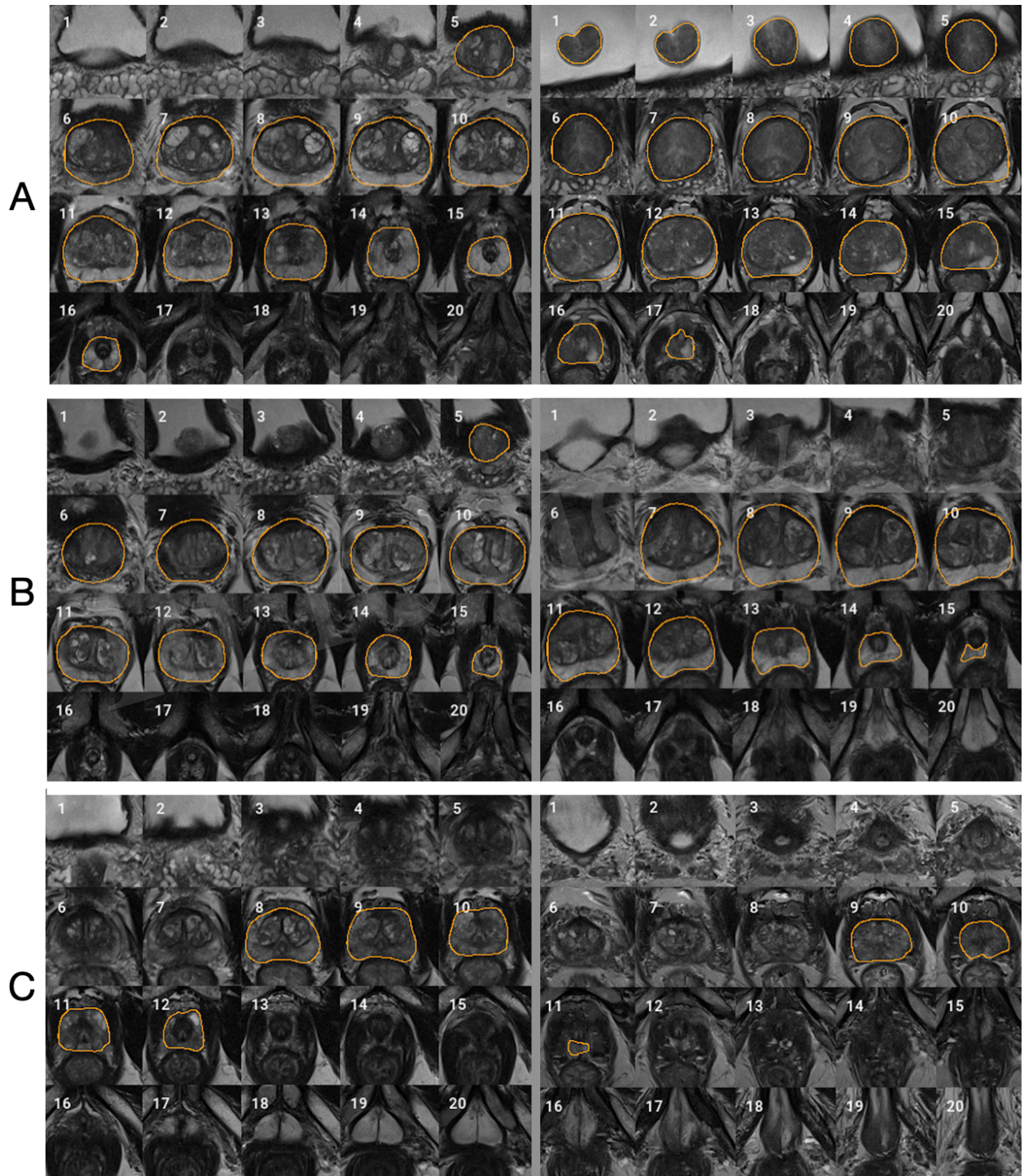
405

406

407

408

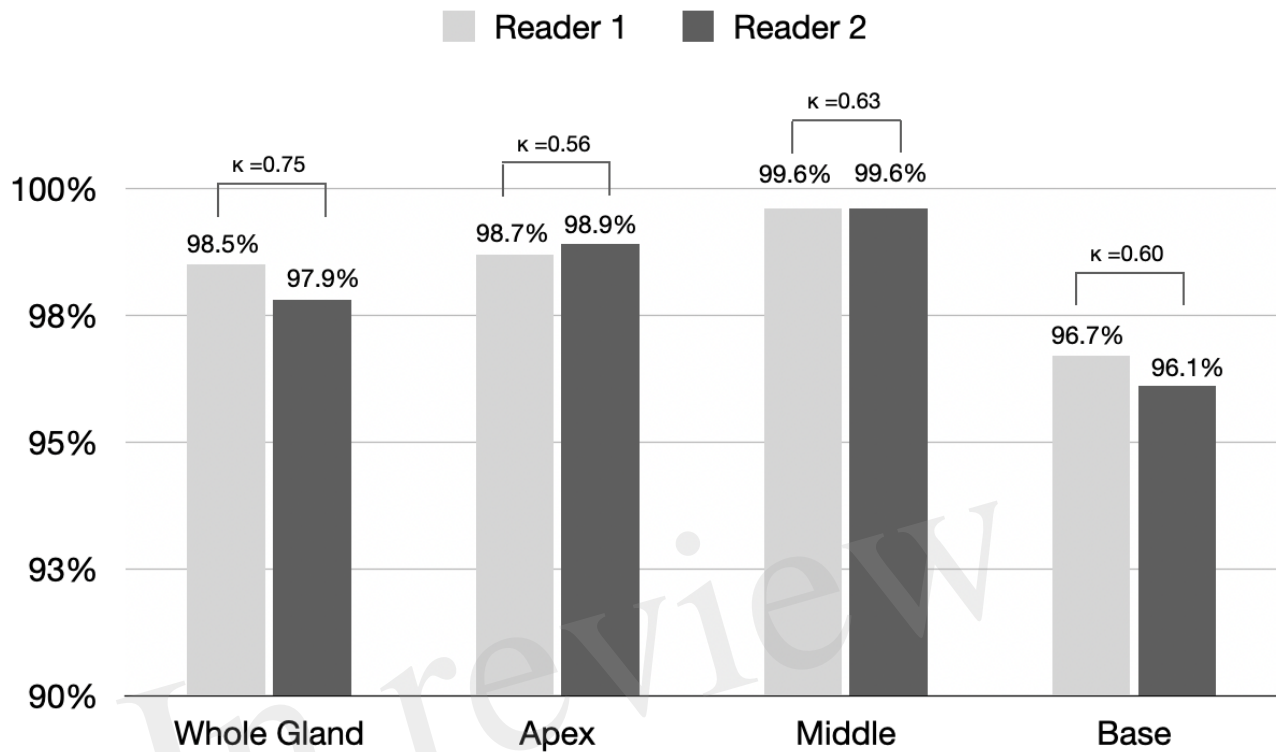
409



410

411 **Figure 2:** Typical examples for each visual grade. Row A, B, and C represent two segmentation
 412 examples with visual grades 3 (excellent), 2 (acceptable), and 1 (unacceptable), respectively. Slice 1-
 413 20 represents MRI slices from superior to inferior. Regions encircled by organ boundary are the
 414 prostate whole gland.

415



416

417 **Figure 3:** The proportion of segmentation with acceptable or excellent performance evaluated by
 418 radiologists 1 and 2 among all MRI scans (n=3210). Kappa statistics between the two readers were
 419 also provided in the figure.

420

421

422

423

424

425

426

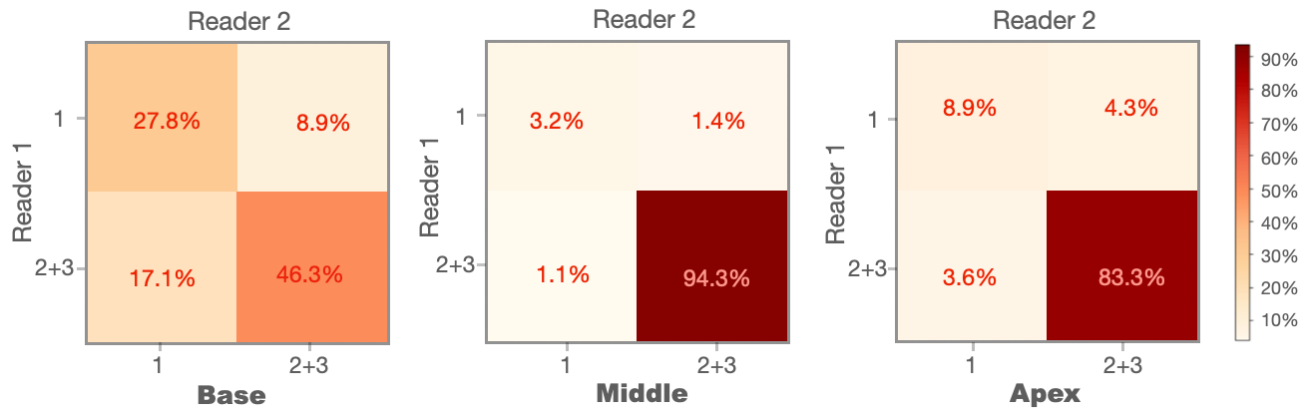
427

428

429

430

431



432

433 **Figure 4:** Confusion matrices of the prostate base, midgland, and apex for the cases without excellent
 434 segmentation (n=281).

435

436

437

438

439

440

441

442

443

444

445

446

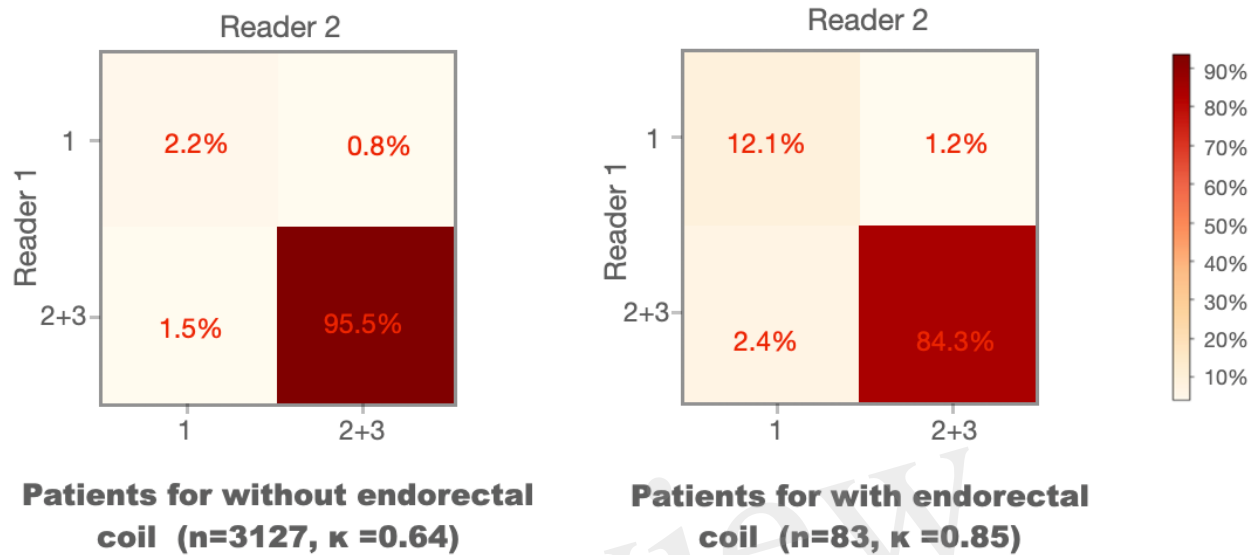
447

448

449

450

451



452

453 **Figure 5:** Confusion matrices of the visual grades of segmentation on MRI scans with and without
 454 endo-rectal coils. Kappa coefficient (κ) is used to measure the inter-rater variability between the two
 455 readers.

456

457

458

459

460

461

462

463

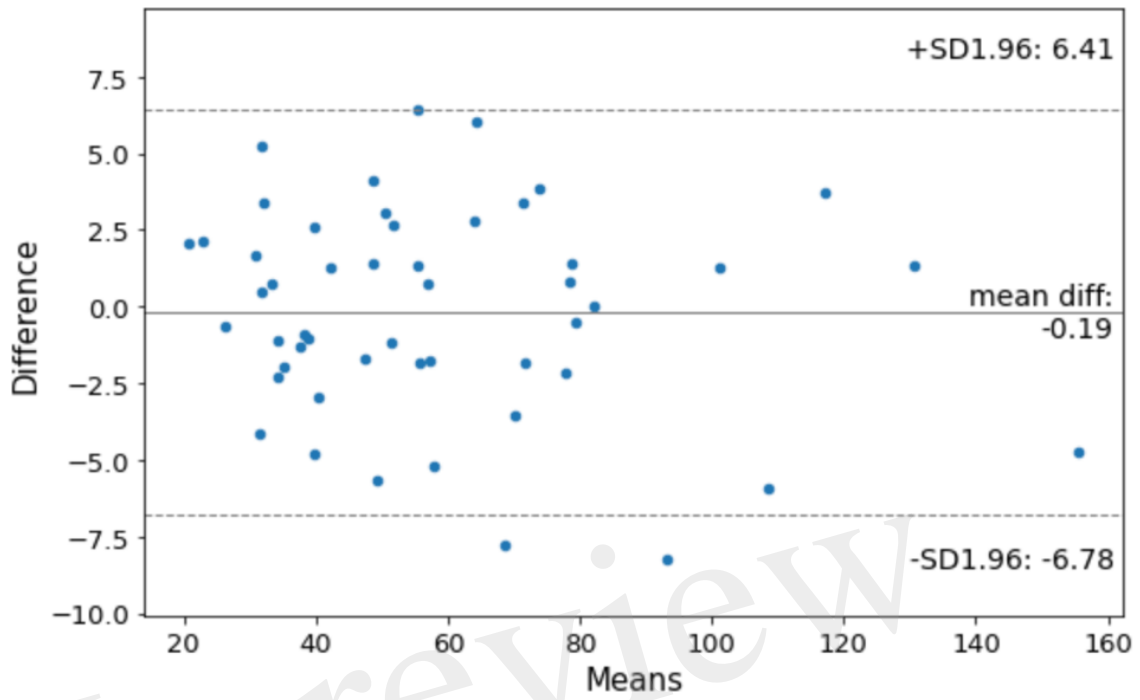
464

465

466

467

468



469

470 **Figure 6:** Bland–Altman plot to show the agreement between manual and DANN-enabled WPG
471 volume measurements.

472

Figure 1.JPEG

In review

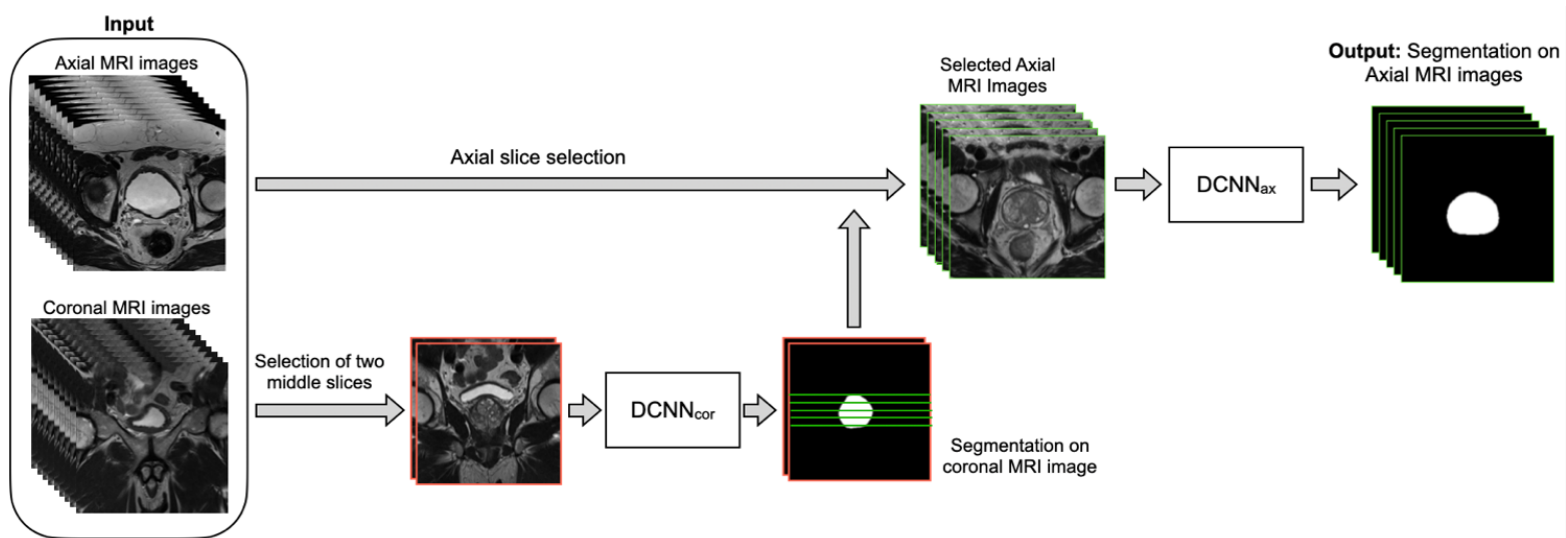


Figure 2.JPEG

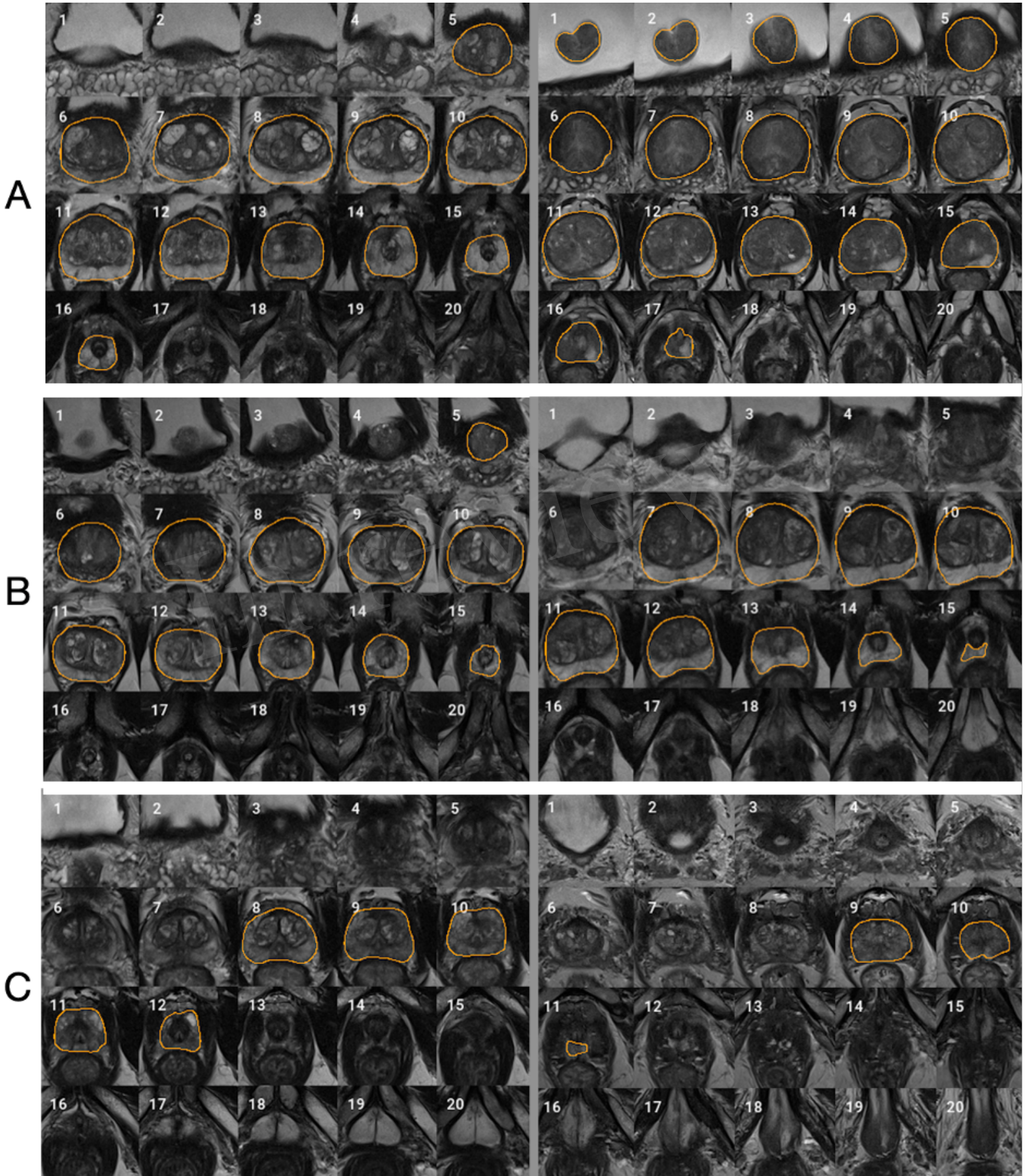


Figure 3.JPEG

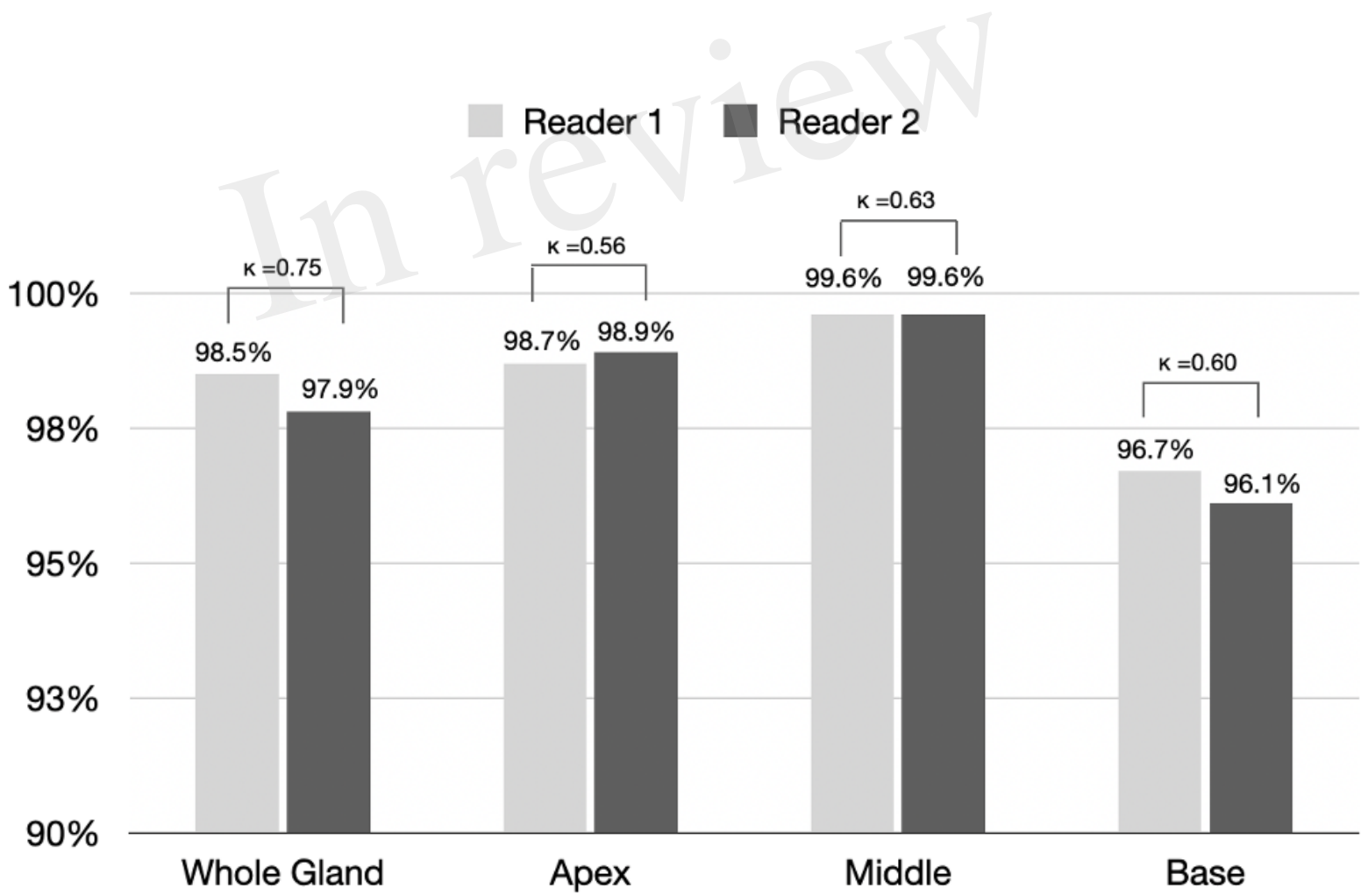
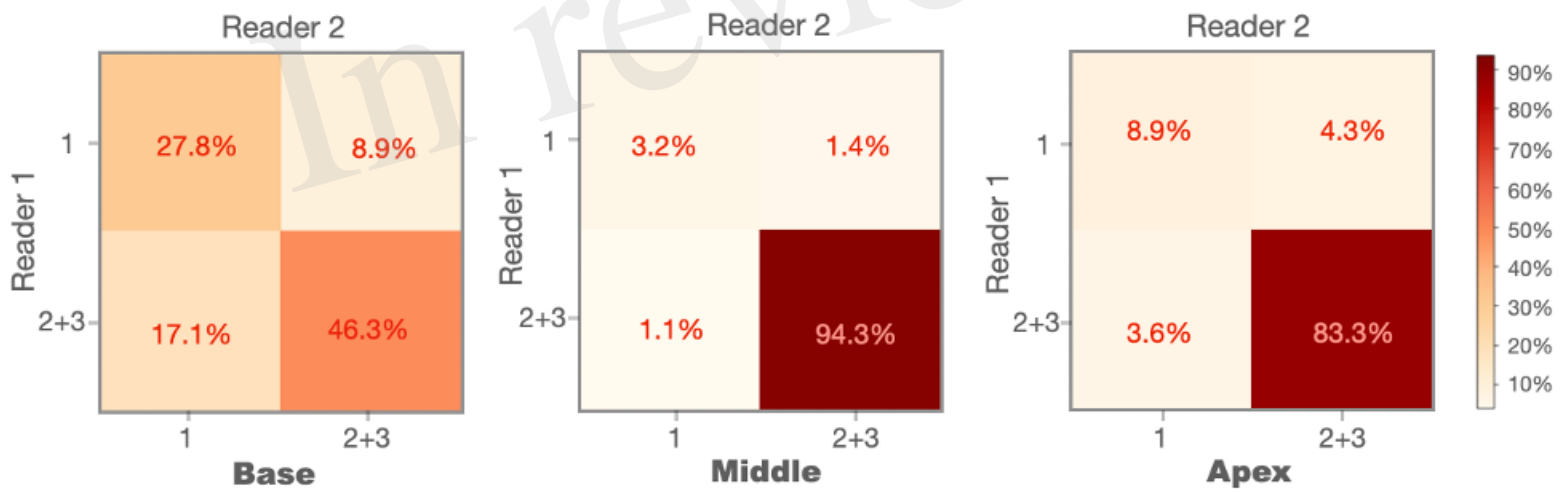


Figure 4.JPEG



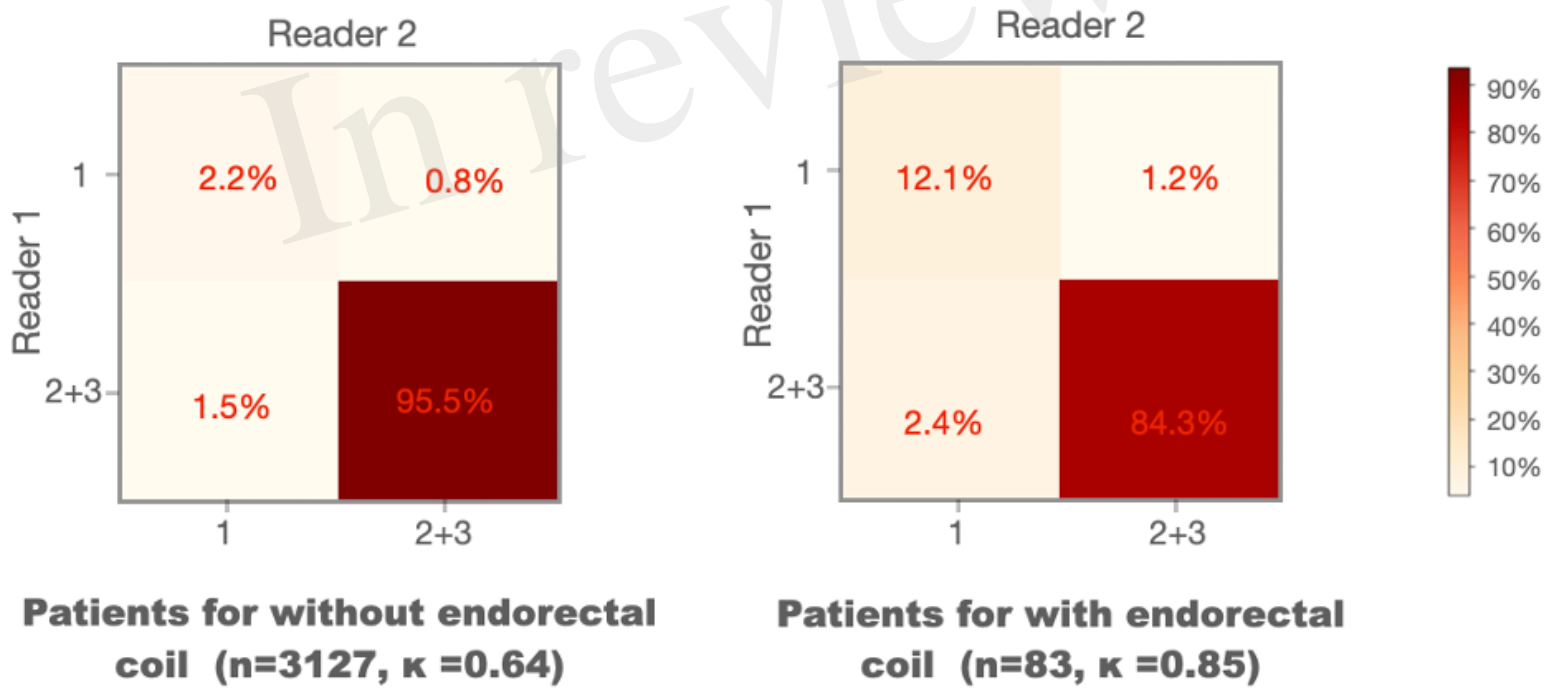


Figure 6.JPEG

In review

