

1
2
3
4
5
6
7
8
9
10
11
12
13
14

A Comprehensive AI Model Development Framework for Consistent Gleason Grading

Supplementary Material

Supplementary Methods2

Supplementary Tables8

Supplementary Figures10

Supplementary Notes.....16

15 **Supplementary Methods**

16 **A!MagQC: A quantitative digital pathology image quality control solution**

17 With the increasing use of digital pathology, vast amounts of data are generated on a daily basis.
18 However, there are common quality issues (as shown in Supplementary **Figure 1a**), and
19 visually assessing image quality has become a tedious and heavy workload for researchers.
20 While perceptual image quality estimators based on perception-based image quality evaluators
21 (PIQE) can calculate a no-reference image quality score, this approach has proven to be less
22 effective for histological images. Although some tools have been developed previously for
23 histological images, they have been limited to evaluating only Haematoxylin and Eosin (H&E)
24 images.

25 To address this gap, A!MagQC was developed to provide fully automated quality
26 control for any histologically relevant imaging modality, including Haematoxylin and Eosin
27 (H&E), Immunohistochemistry (IHC), and Multiplexed Fluorescence (MF). The software
28 automatically detects the image size (magnification) and type from the metadata of each image
29 file. The user interface is shown in **Supplementary Figure 1b**.

30 The first step in assessing the quality of tissue images is to detect the tissue and separate
31 it from the background. This is achieved by applying adaptive thresholding. To evaluate the
32 quality of the tissue at a local level, we performed a parallel analysis of tiles measuring 256*256
33 pixels throughout the Region of Interest (ROI). We have identified five relevant features to
34 differentiate local quality in whole slide images, as shown in **Supplementary Figure 1c**:

- 35 • Focus: We quantify the focus in an image using the Variance of the Laplacian
36 Transform. The Laplacian operator measures the second derivative of an image,
37 highlighting regions of an image with sharp intensity changes. High variance of
38 intensity change, which indicates sharp and smooth changes, is representative of a

39 normal, in-focus image. Conversely, low variance indicates an image with few sharp
40 edges, typically an out-of-focus image.

41 • Contrast: We quantify Contrast by measuring the difference between the top 1% of
42 high-value positive pixels and the bottom 1% of low-value positive pixels in each tile.
43 The range must be high enough for good separation of nuclei signal and background.

44 • Saturation: We measure the percentage of pixels that have a maximum intensity value,
45 which is 255 for 8-bit unsigned integer digital image pixels.

46 • Artifacts: The main structure of interest in Histology images is usually the nuclei. The
47 morphological open operation, which is an erosion followed by a dilation with the same
48 structuring element for both operations, is used to perform an image opening. If the
49 structuring element or kernel is bigger than the average nuclei size, it highlights dirt
50 and blurry objects that sometimes occur in the images.

51 • Texture Uniformity: Computing the Gray Level Co-occurrence Matrix (GLCM)
52 calculates how often a pixel with gray-level (grayscale intensity) value "i" occurs
53 horizontally adjacent to a pixel with the value "j". Measuring the uniformity of the
54 pixels allows us to highlight regions with different densities of nuclei. Notably, visceral
55 fat tissue surrounding organs of interest often has a very different texture.

56 The pipeline was designed for multiplexed fluorescence images, with the algorithm
57 directly applied to the grayscale image of the DAPI fluorescence signal. H&E images are
58 converted to optical density (OD) using a logarithmic transformation before analysis.

59

60 **A!HistoClouds: The cloud-based digital pathology image annotation and management**
61 **platform**

62 AI-driven computational pathology diagnosis is an emerging but rapidly developing field. It
63 uses computational algorithms to classify cancer and other diseases, based on the annotated
64 images. Annotating pathological images requires experienced pathologists with years of
65 training. A high-quality annotated image database is the basis for developing AI-based
66 diagnostic solutions, because most successful models are derived from supervised learning. It
67 is important to mark and annotate specific areas/structures/features to describe a disease at the
68 cellular level, and then build and validate the models. Currently, there is no effective "medium"
69 to transfer a pathologist's knowledge and experience to a machine. A!HistoClouds is a cloud-
70 based structural annotation platform designed to enable pathologists to address this gap.

71 The image viewer (See **Supplementary Figure 2a**) based on the openseadragon
72 software library can visualize DP images with high resolution of 40x objective lens, load image
73 blocks quickly and smoothly, without consuming a lot of device memory and Internet data.
74 The most important basic event functions, such as "zoom", "pan" and "home page", can all be
75 customized using its application programming interface (API) as illustrated in See
76 **Supplementary Figure 2b-c**.

77 Annotation tool is one of the basic components of A!HistoClouds. The annotation tool
78 provides a ROI management system and has the ability to create an adjustable ROI on top of
79 the image viewer. In other words, when the entire image is moved by the user, the ROI adheres
80 to a specific area on the image. Once the ROI shape is released, it can be fine-tuned. ROI
81 management system refers to a way to easily manipulate and manage many ROIs on the viewer.

82 A!HistoClouds provides three ROI drawing methods for annotating tasks in the viewer,
83 which can be found in the toolbox button as illustrated in **Supplementary Figure 2d**. They are
84 the freehand drawing, polygonal-dot drawing and brush drawing shown in **Supplementary**
85 **Figure 2e**. When selecting ROI for further operation, user can click the right-click menu panel.
86 They include labelling, copying, attribute updating and deleting operation. When user click

87 "More label", ROI can be renamed (default label is "unknown"), and a tag window dialog box
88 will appear for naming choices, as shown in **Supplementary Figure 2f**. Multiple ROI selection
89 is one of the great features love to be used by our pathologist (See **Supplementary Figure 2g**).
90 They first draw many ROIs and then label them at once, which is very helpful for them to save
91 a lot of valuable annotation time.

92 When pointing the mouse at them, user can easily identify the ROI and related
93 information on the ROI panel as demonstrated in **Supplementary Figure 2H**. In addition,
94 selected ROI can be hidden and shown easily by click on the "eye" icon at the ROI panel as
95 shown in **Supplementary Figure 2i**. This is a particularly useful feature that allows
96 pathologists to draw ROIs of tissues that may be obscured by another large ROI.

97 For AI-assisted diagnosis and semi-automatic annotation, the outputs generated by the
98 AI model can be converted into ROI in A!HistoClouds. Therefore, Pathologists can view and
99 modify the ROI annotations using the A!HistoClouds image viewer, and fine-tune accordingly.

100 Besides, the time spent on annotation is an important measure to understand how easy
101 it is for the pathologist to annotate the entire image slice. They are evidence showing the time
102 spent by pathologists on fully manual annotations and time spent on some fine-tuning (semi-
103 automatic annotations) of the ROI generated based on the AI model. Therefore, A!HistoClouds
104 will automatically record the time spent to draw in each ROI (See **Supplementary Figure 2b**)
105 for performance evaluation.

106

107 **Hardware and software for model development**

108 We performed AI model training and testing on MATLAB 2021a (MathWorks Inc., USA) with
109 its Deep Learning and Deep Learning, Image Processing and Parallel Computing toolboxes on
110 the Windows 10 X64 operating system. The computer specifications are RAM: 1.0TB, CPU:

111 Intel(R) Xeon(R) Gold 6242 CPU @ 2.80GHz, and GPU: single NVIDIA Tesla V100-PCIE-
112 16GB.

113

114 **AI model selection and optimization**

115 Several models of different architectures were initially trained using high-resolution patches
116 (1.12 μ m/pixel). **Supplementary Figure 3a-c** indicates that NasNet Mobile (macro F1 = 0.68)
117 is slightly inferior to ResNet50 (macro F1 = 0.71) and Vgg16 (macro F1 = 0.71) in terms of
118 macro F1 score, with the differences between ResNet50 and Vgg16 being subtle. Given that
119 ResNet50 has less parameters, it was selected for faster deployment. The models based on
120 ResNet50 structures were trained at four different scales and then applied to test images to
121 compare their performance and optimize the scale factors.

122

123 **Evaluation of AI model on annotation-level and WSI-level using multiple pathologists’ 124 annotations as reference**

125 In this study, the ground truth annotations were reviewed and adjusted by multiple pathologists
126 based on NUH annotations, resulting in different set of annotations, as shown in
127 **Supplementary Figure 4a and b**. Annotations agreed upon senior pathologists were used to
128 assess the models’ performances and select the optimal model. Besides, we evaluated the AI
129 model performance using NUH and 9 pathologists’ annotations respectively (**Supplementary
130 Figure 4c**). Despite these variations, the model demonstrated superb performance in
131 identifying non-malignant tissues. Inter-observer variations not only exist on annotation level,
132 but also Gleason grading on WSI-level. **Supplementary Figure 4d** demonstrated the Grade
133 Groups (GGs) determined by different pathologists. We assessed the consistency of GG among
134 different pathologists and AI model using Quadratic Weighted Kappa, shown in
135 **Supplementary Figure 4e**.

136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160

Three-phase clinical validation of AI-assisted diagnosis

Although the AI model performs well on the image data set, it is imperative to conduct further validation to assess its practical utility in assisting pathologists in real-world application. In our study, we designed a comprehensive three-phase experiment with the objective of comparing the efficacy and efficiency of Gleason Grading through microscopic examination, whole slide image (WSI) examination with and without AI assistance.

For this experiment, we randomly selected 19 slides from the test set, ensuring a representative sample. To establish a reliable ground truth, the Gleason Grade Groups for these slides were independently determined by four senior pathologists. These senior pathologists' assessments were utilized as the reference for calculating the Quadratic Weighted Kappa.

In each phase of the experiment, pathologists meticulously examined the 19 selected slides individually. They assessed and assigned Gleason Scores to each slide while recording the time spent on evaluation. The WSIs were captured at 20× magnification (0.5 μm/pixel) using Akoya Biosciences Vectra Polaris scanner. Phase 3 introduced AI assistance, which encompassed a range of features, including pseudo annotation, tumor percentage, Gleason Pattern percentage, and Gleason Score, all generated by our AI model.

In phase 1, only three pathologists from Singapore participated due to the logistical challenge of shipping glass slides to China. To limit recall bias, we ensured that for each phase and for each pathologist, the order of slide review was intentionally randomized. We also provided comprehensive user guides and pre-experiment training to ensure that all participants were proficient in using A!HistoClouds. To maintain methodological integrity, we implemented a mandatory washout period of at least 20 days between each phase. Additionally, in phases 2 and 3, the filenames of the whole slide images (WSI) were randomly generated, respectively.

161 **Supplementary Tables**

162 **Supplementary Table 1 Patient Characteristic** Profile of 214 patients included in the study. One
 163 patient's information is missing. Both prostatectomy specimens and biopsy samples were collected
 164 from 103 patients, and the other patients provided either prostatectomy specimen or biopsy sample.
 165

Age	Number	Percentage
45–50	1	0.5%
51–60	23	10.7%
61–70	132	61.7%
71–80	49	22.9%
81-90	9	4.2%
Gleason Score		
3+3	13	6.1%
3+4	82	38.3%
4+3	58	27.1%
4+4	4	1.9%
3+5	5	2.3%
5+3	2	0.9%
4+5	35	16.4%
5+4	8	3.7%
5+5	7	3.3%

166

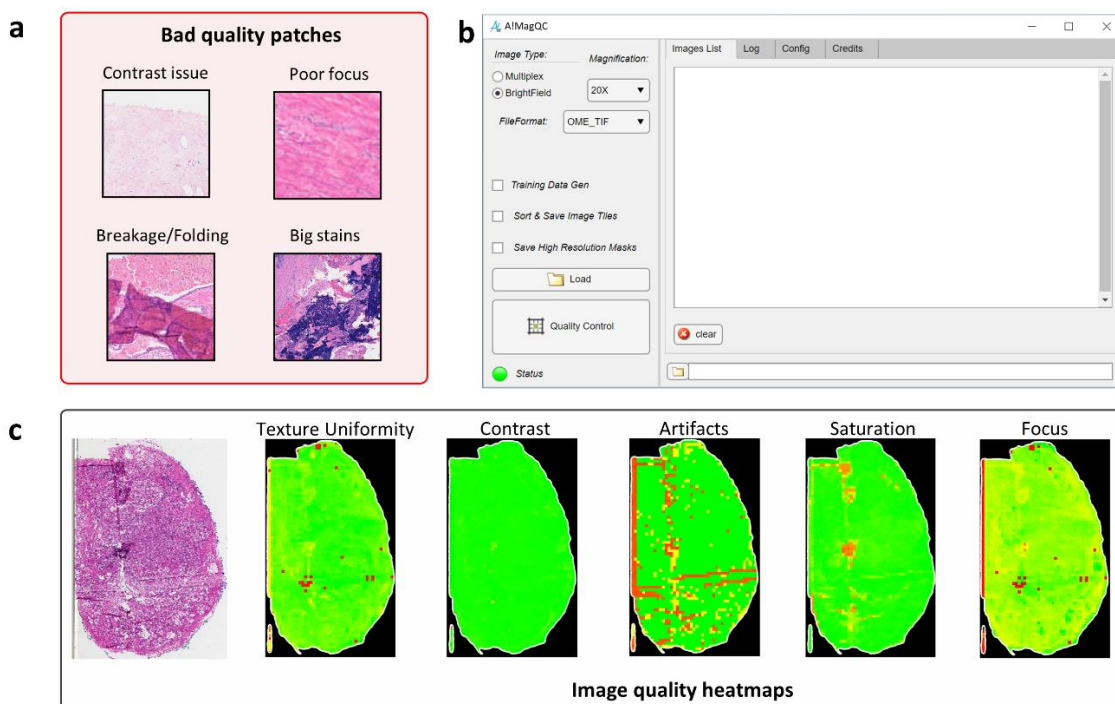
167 **Supplementary Table 2 Configuration of color augmentation** All values are subject to $\pm 5\%$
 168 variance. R, G, B values are first adjusted by addition/subtraction, then rescaled to [0 1],
 169 followed by clipping, in which values below Low_in are mapped to 0 and values
 170 above High_in map to 1. Low_in and High_in values apply to all R, G, B channels.

Configuration	R value	G value	B value	Low_in	High_in
1	-60	-50	-20	0.05	0.95
2	-30	-45	-60	0.05	0.95
3	+35	+70	+35	0.1	0.9

171

172 **Supplementary Figures**

173



174

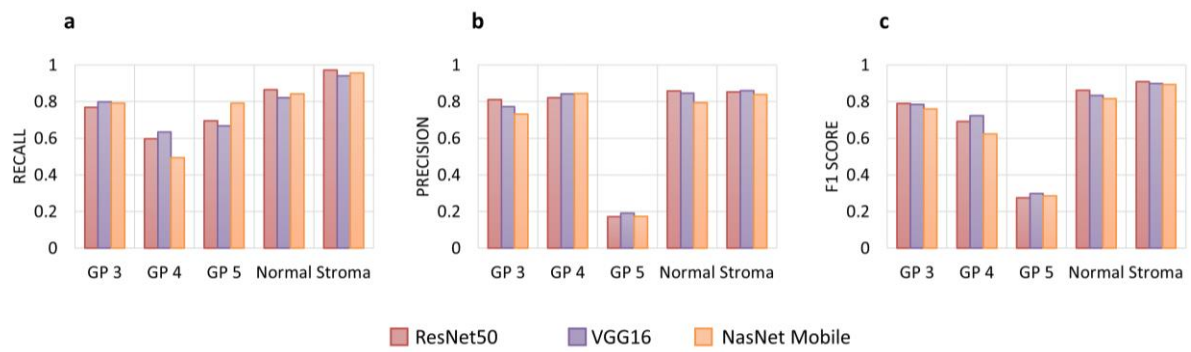
175 **Supplementary Figure 1 Overview of A!MagQC** (a) Some examples of common quality
176 issues of histopathological images. (b) User interface of A!MagQC. (c) Heatmap generated by
177 A!MagQC that identify different quality issues. User can easily locate and check the low-
178 quality patches according to the heatmap.

179
180



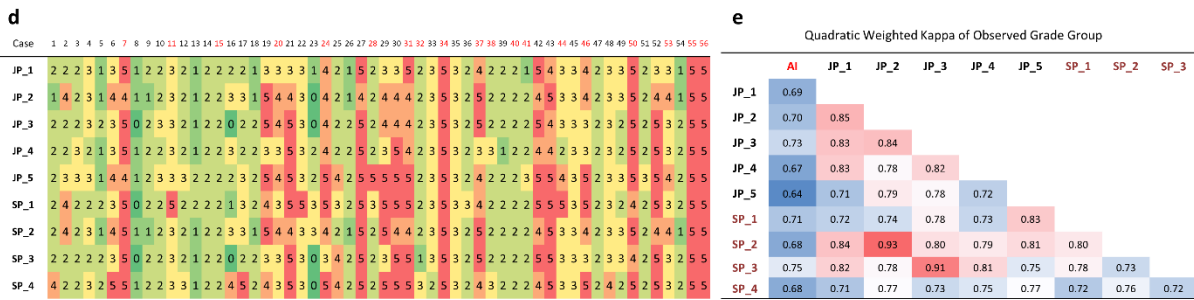
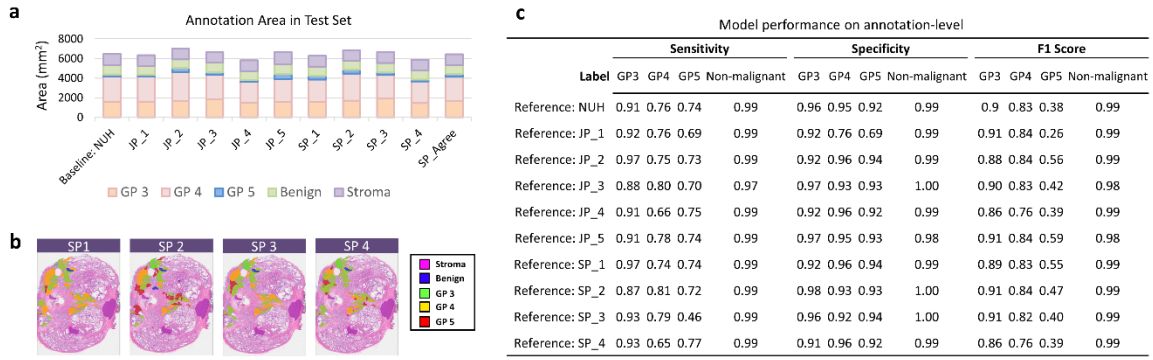
181

182 **Supplementary Figure 2 Overview of A!HistoClouds** A!HistoClouds consists of: (a) image
183 viewer. (b)-(c) timer and basic event function, such as "zoom", "pan" and "home page". (d)
184 toolbox of ROI drawing tools (e) drawing tools: freehand drawing, polygonal-dot drawing and
185 brush drawing. (f) label selection window, (g) annotation panel. (h)-(i) panel of each ROI,
186 where user can easily find the related information of the ROI, and hide the ROI by clicking on
187 the "eye" icon at the ROI panel.

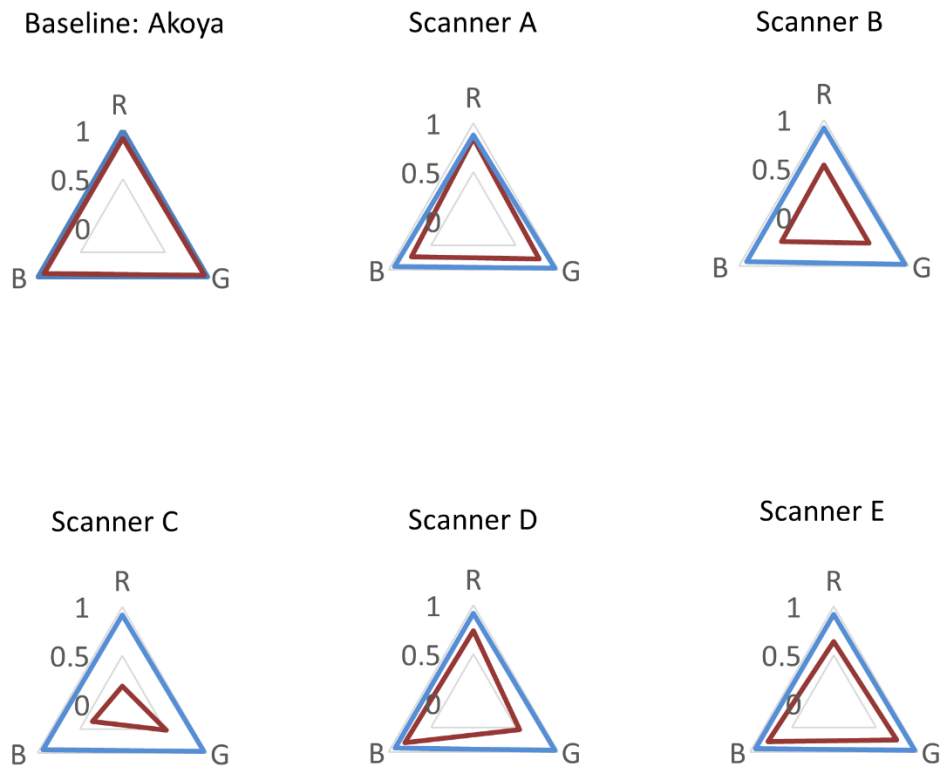


189

190 **Supplementary Figure 3 Model Selection** To select the network architecture, we used high-
 191 resolution (1.12 μ m/px) image patches to train three different models separately and compared
 192 their performances. Patch-level performances on test set were shown in (a)-(c). ResNet50 was
 193 selected as the preferred model due to its superior performance and smaller network size.



196 **Supplementary Figure 4 Evaluation of Model Performance on Prostatectomy Specimens**
 197 **Using Multiple Pathologists' Annotations** The AI model was tested on prostatectomy
 198 specimens and compared with annotations made by multiple pathologists to evaluate its
 199 performance on both annotation- and WSI-level. The number of annotations made by different
 200 pathologists and the agreed upon annotations are presented in (a). Inconsistent annotations
 201 made by different pathologists are illustrated in (b), leading to variations in sensitivity,
 202 specificity, and F1 score when different standards were applied, as shown in (c). Despite these
 203 variations, the model demonstrated superb performance in identifying non-malignant tissues.
 204 On the WSI level, GGs determined by different pathologists were summarized in (d). The
 205 model achieved a weighted kappa of 0.71 on average with four senior pathologists, while the
 206 average weighted kappa among the four pathologists was 0.75, as shown in (e).



207

— Before image appearance migration — After image appearance migration

208

Supplementary Figure 5 Histogram intersection between baseline dataset and other scanner datasets before and after image appearance migration

209

To quantify the effect of image appearance migration, we measured the image similarity before and after migration

210

using histogram intersection of R, G and B channel between baseline and the others. The results

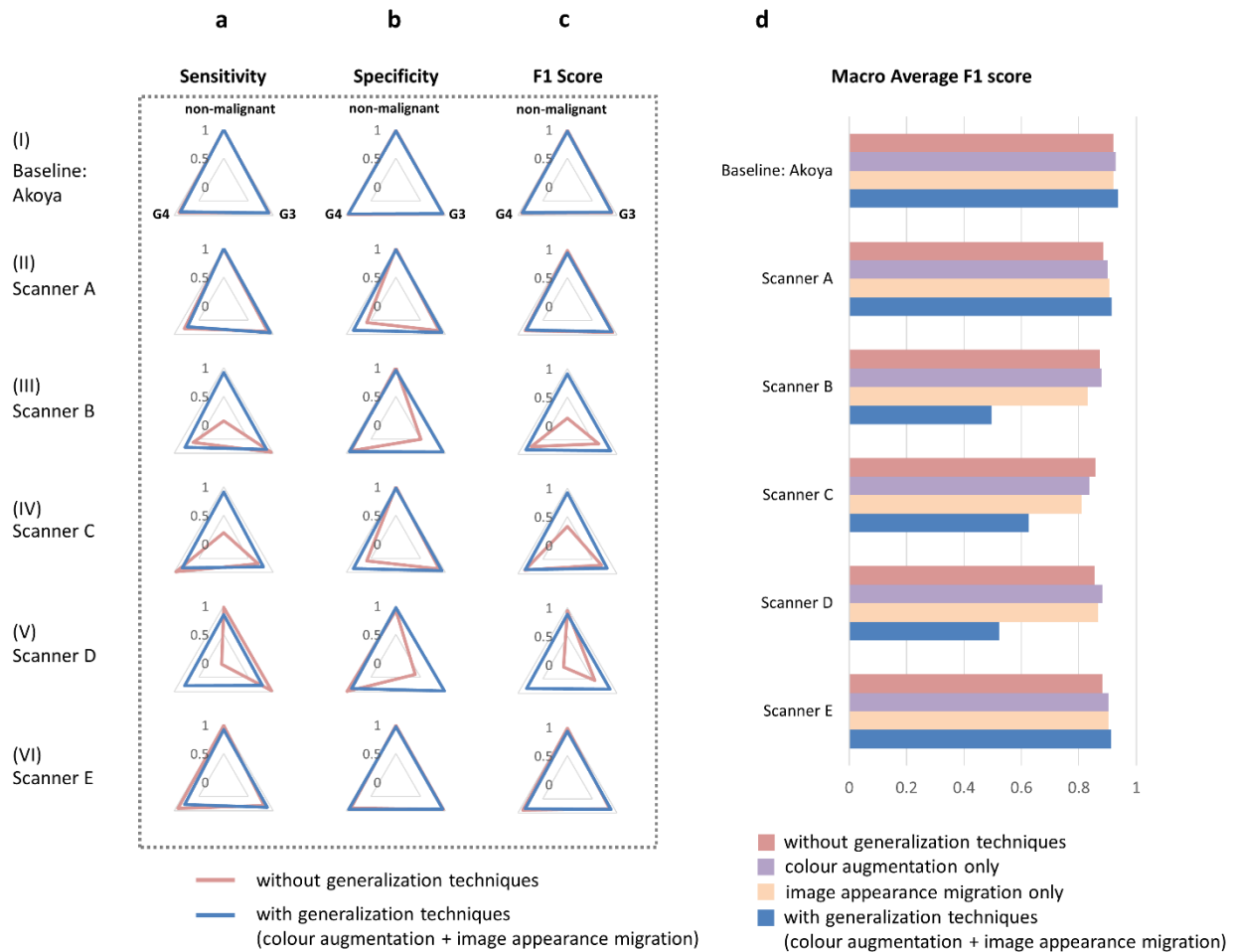
211

showed that migration increased the similarity between baseline and images acquired from

212

other scanners, with almost perfect overlap in histogram intersection for all channels.

213



215
216

217 **Supplementary Figure 6 Details of model performance across different scanners (a)-(c)**
218 Statistical evaluation shows significant improvements in sensitivity, specificity, and F1 score
219 across all classes after applying generalization techniques. (d) The macro average F1 score of
220 each scanner dataset across various generalization techniques is presented. The generalization
221 techniques implemented in this study consist of color augmentation and image appearance
222 migration, and their effects were assessed separately.

223 **Supplementary Notes**

224 **Supplementary Note 1 Pseudo code of train-test splitting** Considering that prostatectomy
 225 specimens are much larger and thus contain more information than biopsies, we used the
 226 annotations made by pathologists from NUH on prostatectomy WSIs to train our models. The
 227 187 radical prostatectomy WSIs were split into training and testing sets, whilst the annotated
 228 biopsy images were used for testing only. The training and testing set split ratio of
 229 prostatectomy WSIs is 7:3 for the number of WSIs, evenly divided to ensure the same ratios of
 230 areas for each annotated class in both training and testing since the annotated area of different
 231 classes may vary significantly from slide to slide.

233 **Algorithm:** Stochastic Search for balanced dataset

234 **Data:**

235 $D = \{(X_i, Y_i)\}_{i=1, \dots, 187} \leftarrow$ dataset with 187 images,
 236 $X_i \leftarrow$ prostatectomy specimens image,
 237 $Y_i \leftarrow$ collection of annotations on patch-level,
 238 $C = \{G3, G4, G5, Stroma, Normal\} \leftarrow$ class label sets

239 **Output:**

240 $D_{train} \leftarrow$ training dataset, 70% of D
 241 $D_{test} \leftarrow$ testing dataset, 30% of D

242
 243 $FOUND = False$

244 **while** not FOUND **do**

245 **Step A:**

246 Random shuffle the D, let

247 $D_{train} \leftarrow \{(X_j, Y_j)\}_{j=1, \dots, 132}$

248 $D_{test} \leftarrow \{(X_k, Y_k)\}_{k=133, \dots, 187}$

249
 250 **Step B:**

251 Calculate the number of each class on patch-level in training
 252 dataset and test dataset respectively, i.e.

253 $N_{train}^{Stroma}, N_{train}^{Normal}, N_{train}^{G3}, N_{train}^{G4}, N_{train}^{G5},$

254 $N_{test}^{Stroma}, N_{test}^{Normal}, N_{test}^{G3}, N_{test}^{G4}, N_{test}^{G5}$

255
 256 **Step C:**

257 Check the ratio D_{train}/D is $70\% \pm 5\%$

258 **for** Class in C **do**

259 $Ratio^{class} = N_{train}^{class} / (N_{train}^{class} + N_{test}^{class})$

260 **if** $65\% \leq Ratio^{class} \leq 75\%$ **then**

261 FOUND = True

262 **else**

263 FOUND = False

264 Break

265 **end**

266 **end**

267
 268
 269
 270

271 **Supplementary Note 2 Pseudo code of voting algorithm** During the testing phase, the trained
 272 model was applied to test images through a sliding window operation. To ensure that the
 273 detection was comprehensive, the window overlap was set at 50%, resulting in the centre box
 274 being shared by four consecutive windows. A voting strategy was subsequently employed to
 275 determine the label and probability score of each centre box. The final label was determined
 276 based on the most frequently occurring label among the four windows. If there was no such
 277 label, the final label was chosen based on its higher probability score. The probability score of
 278 the center box was then computed as the mean score of the windows corresponding to the final
 279 label.

280

281 **Algorithm:** Voting algorithm for each overlapped patch

282 **Input:** $D = \{(C_i, S_i)_{i=1, \dots, 4} : C_i \in \{G3, G4, G5, Stroma, Normal\},$
 283 $S_i \in (0, 1),$ a collection of label C_i and score S_i for each overlapping
 284 patch.
 285 **Output:** The label C and score S for each overlapped patch.
 286 **If** $|Set \{C_1, C_2, C_3, C_4\}| = 4,$ **then**
 287 | Class: $C = \underset{C}{\operatorname{argmax}} \{S_i : \{C_i, S_i\}\}$
 288 | Score: $S = S_i$
 289 **else if** Most frequent class $C,$ **then**
 290 | Class: $C = C$
 291 | Score: $S = \operatorname{mean}(S \mid \text{most frequent } C)$
 292 **else**
 293 | Class: $C = \underset{C}{\operatorname{argmax}} \{E((S \mid C_i))\}$
 294 | Score: $S = \operatorname{mean}(S \mid C)$

295