

1 **Pixelwise H-score: a novel digital image analysis-based metric to quantify**
2 **membrane biomarker expression from immunohistochemistry images**

3 **Sripad Ram^{1*}, Pamela Vizcarra², Pamela Whalen², Shibing Deng³, CL Painter², Amy Jackson-**
4 **Fisher², Steven Pirie-Shepherd², Xiaoling Xia^{2,4} and Eric L. Powell²**

5 ¹Drug-Safety Research and Development, Pfizer Inc., La Jolla, CA 92121.

6 ²Tumor Morphology Group, Oncology Research and Development, Pfizer Inc., La Jolla, CA 92121.

7 ³Biostatistics Unit, Oncology Research and Development, Pfizer Inc., La Jolla, CA 92121.

8 ⁴Current affiliation: Ventana Medical Systems, Tucson, AZ 85755.

9

10 ***Corresponding author:** Pfizer Inc., 10646 Science Center Drive, San Diego, CA 92121. Email:

11 Sripad.ram@pfizer.com; Phone: 1-858-622-5904; Fax: 1-858-678-8124

12

13 **Running title:** Pix H-score algorithm

14

15

16 **ABSTRACT**

17 Immunohistochemistry (IHC) assays play a central role in evaluating biomarker expression in tissue
18 sections for diagnostic and research applications. Manual scoring of IHC images, which is the current
19 standard of practice, is known to have several shortcomings in terms of reproducibility and scalability to
20 large scale studies. Here, by using a digital image analysis-based approach, we introduce a new metric
21 called the pixelwise H-score (pix H-score) that quantifies biomarker expression from whole-slide scanned
22 IHC images. The pix H-score is an unsupervised algorithm that only requires the specification of intensity
23 thresholds for the biomarker and the nuclear-counterstain channels. We present the detailed
24 implementation of the pix H-score in two different whole-slide image analysis software packages
25 Visiopharm and HALO. We consider three biomarkers P-cadherin, PD-L1, and 5T4, and show how the pix
26 H-score exhibits tight concordance to multiple orthogonal measurements of biomarker abundance such
27 as the biomarker mRNA transcript and the pathologist H-score. We also compare the pix H-score to
28 existing automated image analysis algorithms and demonstrate that the pix H-score provides either
29 comparable or significantly better performance over these methodologies. We also present results of an
30 empirical resampling approach to assess the performance of the pix H-score in estimating biomarker
31 abundance from select regions within the tumor tissue relative to the whole tumor resection. We
32 anticipate that the new metric will be broadly applicable to quantify biomarker expression from a wide
33 variety of IHC images. Moreover, these results underscore the benefit of digital image analysis-based
34 approaches which offer an objective, reproducible, and highly scalable strategy to quantitatively analyze
35 IHC images.

36

37

38

39 INTRODUCTION

40 Immunohistochemistry (IHC) is a core technology that is used to evaluate the spatial distribution and
41 abundance of biomarkers at the protein level in tissue samples. In oncology clinical diagnosis and research
42 applications, IHC assays play a central role in tumor characterization and biomarker assessment. Typically,
43 IHC images are qualitatively evaluated by a trained expert, such as a pathologist, and in some cases this is
44 complemented by a semi-quantitative score [1]. However, visual quantitative scoring of IHC images is not
45 routinely performed due to several shortcomings. On the one hand, visual quantitative scoring is time
46 consuming and is often not feasible to perform on a routine basis especially for large studies. On the other
47 hand, visual quantitative scores are subjective and often have a limited dynamic range due to their
48 categorical nature (e.g. manual scores of 0, 1+, 2+, and 3+). Consequently, they may not have the
49 granularity to adequately capture biomarker expression from an IHC slide [2, 3]. The subjectivity of the
50 scoring process, in turn, can manifest as poor inter- and intra-observer concordance, and this has been
51 the subject of numerous studies [4-8]. While concordance in visual quantitative scoring can be improved
52 by the development of standardized scoring guidelines and extensive training [9, 10], the labor-intensive
53 aspect and the limited dynamic range still remain as major impediments to the widespread use of visual
54 quantitative scoring of IHC images.

55 Digital image analysis (DIA) based tools overcome some of these limitations of visual quantitative scoring
56 by enabling fast, objective, and highly reproducible quantification of biomarkers from whole-slide IHC
57 images [1, 11]. DIA endpoints are typically continuous variables (e.g. cell density and % positive cells) and
58 offer adequate dynamic range to represent biomarker expression in the IHC image. One of the widely
59 used endpoints to quantify biomarker expression is the H-score [2, 12]. In the H-score algorithm (Figure
60 1A) individual cells and their sub-cellular compartments (i.e. nucleus, cytoplasm, and cell membrane) are
61 first detected, and based on the relative expression of the biomarker of interest in one or more sub-

62 cellular compartments the cells are classified as either positive or negative. The positive cells are further
63 classified into high (3+), medium (2+), or low (1+) based on the biomarker signal intensity. The H-score is
64 given by the ratio of the weighted sum of the number of positive cells to the total number of detected
65 cells. The H-score captures both the intensity and the proportion of the biomarker of interest from the
66 IHC image and comprises values between 0 and 300, thereby offering a dynamic range to quantify
67 biomarker abundance. A different scoring method developed to quantify estrogen and progesterone
68 receptors in breast cancers, the Allred score [2, 12], assigns separate categorical scores for the intensity
69 (0-3) and the proportion (0-5) of the biomarkers in immunolabeled specimens, and the final score is the
70 sum of these two scores. Compared to the H-score, the Allred score has a limited dynamic range (0-8) and
71 is not extensively used for purposes other than ER/PR quantification in breast cancer. From a digital image
72 analysis standpoint, both the H-score and the Allred score require the detection of individual cells, and
73 this requires robust nucleus and cell segmentation algorithms for individual nucleus detection and
74 delineation of individual cell boundaries.

75 Another scoring methodology, the average threshold method (ATM), adopts a pixelwise approach for
76 quantifying biomarker abundance [13]. The ATM score does not require the detection of individual nuclei
77 or cells and is solely based on the pixel intensities of the DAB chromogen in the spectrally deconvolved
78 image. Consequently, the calculation of the ATM score is relatively straightforward but at the expense of
79 decreased dynamic range as compared to the H-score.

80 The AQUA score [14] also makes use of a pixelwise strategy for quantifying biomarker expression. Here,
81 the tissue is fluorescently labeled for the biomarker of interest along with a nuclear stain and a cell
82 membrane marker. This in turn allows the generation of pixel masks pertaining to different subcellular
83 compartments (e.g., cell membrane, nucleus, or cytosolic mask). The AQUA score is then calculated by
84 taking the total fluorescence signal of the biomarker of interest for a given subcellular mask (e.g. the cell-

85 membrane mask) and normalizing it by the total area of the mask [14]. The advantage of the AQUA score
86 is that it offers a broad dynamic range. However, the calculation of the AQUA score requires the
87 development of a fluorescence-based multiplex assay which can be time consuming and technically
88 challenging. Moreover, the use of fluorescence readout masks anatomic and morphological information
89 (e.g. necrotic regions, stroma, etc.) that are readily detectable from a brightfield IHC image.

90 In this manuscript, three different scoring methods are compared, which are illustrated in Figure 1. We
91 introduce a new DIA method, the pixelwise H-score (pix H-score), for quantifying biomarker abundance
92 from brightfield IHC images by making use of individual pixel intensities in DAB and hematoxylin channels
93 and leveraging weighted intensity averages. Our motivation behind developing the pix H-score is to create
94 a simple, yet robust metric to accurately quantify biomarker expression without relying on the detection
95 and delineation of individual cells and their sub-cellular compartments. The latter makes the
96 implementation of the pix H-score to be relatively straightforward. The pix H-score can be thought of as
97 an equivalent of the traditional H-score that is applied to pixels rather than to cells. The pix H-score takes
98 values between 0 and 300 thereby providing a dynamic range similar to that of the H-score.

99 We evaluated the performance of pix H-score using IHC images of three different membrane biomarkers
100 P-cadherin, PD-L1, and 5T4. For comparison, we also calculated the ATM score and the DIA H-score for
101 these images, where the latter is a DIA implementation of the traditional H-score. Using the pathologist
102 H-score and biomarker mRNA transcript level (measured using qRT-PCR or NanoString analysis of mRNA
103 in adjacent serial sections) as orthogonal measurements of biomarker abundance, we demonstrate that
104 the pix H-score is either comparable or superior to other DIA endpoints in quantifying biomarker
105 abundance in IHC images. We present the detailed implementation of the pix H-score in two commercial,
106 whole-slide, image analysis software packages, Visiopharm and HALO. We also present an empirical
107 resampling approach to quantitatively assess the ability of the pix H-score to estimate biomarker

108 abundance when it is calculated from select regions within the tumor resection when compared to the
109 whole slide pix H-score. We anticipate that the new metric will have broad applicability and pave the way
110 towards establishing an objective, reproducible strategy to quantify biomarker abundance in IHC images.

111 MATERIALS AND METHODS

112 Previously-developed IHC assays for P-cadherin, PD-L1, and 5T4 were used to immunolabel three cohorts
113 of human tumors. Serial sections from these cohorts were also evaluated for target mRNA via NanoString
114 (P-cadherin and PD-L1) or qRT-PCR (5T4). Following H-scoring of the immunolabeled tumor sections by a
115 pathologist, the concordance between the H-score and mRNA values was evaluated by Spearman
116 correlation. To automate the scoring process through digital image analysis, we implemented several DIA
117 strategies using different software tools. Specifically, we implemented digital H-scoring using QuPath and
118 HALO software packages, the ATM score using Visiopharm software, and the pix H-score, the new digital
119 scoring method, using HALO and Visiopharm software packages. To assess the performance of the various
120 DIA algorithms, we calculated the Spearman's correlation coefficient between each DIA endpoint and two
121 different measurements of biomarker abundance, i.e. the pathologist H-score and the target transcript
122 level as assessed using either NanoString technology or qRT-PCR.

123 **Immunohistochemistry:** All tumor samples used in this study were anonymized specimens from
124 commercial and academic sources that collected the specimens with donor consent under Institutional
125 Review Board-approved procedures. For PD-L1, we used twenty-four cases of routinely collected non-
126 small cell lung carcinoma surgical resections. The SP142 clone of anti-PD-L1 antibody was used as per the
127 manufacturer-recommended protocol. For P-cadherin, we used thirty cases of routinely collected head
128 and neck tumor resections. The P-cadherin IHC assay was developed and optimized on the Dako
129 Autostainer system using a custom anti-P-cadherin antibody that was generated as an analyte specific
130 reagent for use in a clinical diagnostic assay. For 5T4, we used twenty-one cases of routinely collected
131 non-small cell lung tumor resections. The development and validation of the 5T4 IHC assay was reported
132 previously [15]. In all three IHC assays hematoxylin was used as the nuclear counterstain and
133 diaminobenzidine (DAB) was the chromogen that was used to detect the biomarker of interest. P-cadherin

134 and PD-L1 slides were scanned using a Leica Aperio AT2 whole-slide scanner at 20x magnification, whereas
135 5T4 slides were scanned using a Hamamatsu Nanozoomer whole-slide scanner at 20x magnification.

136 **NanoString assay:** Messenger RNA (mRNA) was isolated from two 4-micron FFPE slide sections using
137 FormaPure[®] nucleic acid isolation kit according to manufacturer's instructions with the addition of a DNA
138 digestion step. NanoString technology was used to measure RNA transcript levels using the nCounter
139 assay according to manufacturer's recommended protocols. Custom nCounter CodeSet containing either
140 the CDH3 probe (for P-cadherin) or the CD274 probe (for PD-L1) was used. One hundred nanograms of
141 total RNA was hybridized to the custom panel for 16 to 20 hours at 65°C. Samples were processed using
142 an automated nCounter sample prep station. Cartridges containing immobilized and aligned reporter
143 complex were subsequently imaged and counted on an nCounter Digital Analyzer set for maximum fields
144 of view. Reporter counts were analyzed and normalized using NanoString nSolver Analysis
145 Software. Briefly, raw counts were multiplied by scaling factors proportional to the sum of counts for
146 spiked in positive control probes to account for individual assay efficiency variation, and to the geometric
147 average of the housekeeping gene probes to account for variability in the mRNA content. FFPE sample
148 sets were normalized to the following housekeeping genes; for P-cadherin: FTL, GAPDH, GUSB, HMBS,
149 HPRT1, OAZ1, PCBP1, PFN1, PPIA, PSAP and TBP; and for PD-L1: AMMECR1L, CNOT10, CNOT4, COG7,
150 DDX50, EDC3, EIF2B4, ERCC3, FCF1, FTL, GPATCH3, GUSB, HDAC3, HPRT1, MTMR14, PPIA, SAP130, TBP,
151 TMUB2, and ZNF143.

152 **qRT-PCR assay:** The qRT-PCR reaction was performed using the TaqMan Probe-Based Gene Expression
153 Analysis and ABI ViiA7 Real-Time PCR Systems (Life Technologies) as described previously [15]. Target
154 gene and endogenous controls were run in quadruplicate for each probe set on prefabricated TaqMan
155 low density array cards. For each tumor sample 1000 ng of cDNA was diluted to 55 uL with nuclease-free
156 water and 55 uL of TaqMan gene expression master mix was added (Life Technologies, cat # 4352042). A

157 total of 100 uL of sample was added to each of the 8 ports on a single card, after which the plate was
158 sealed and centrifuged two times in Sorvall/Heraeus buckets based on manufacturer's directions. TaqMan
159 array cards were then sealed and loaded into the ABI ViiA7 thermal cycler and run. Default thermal cycling
160 conditions were as follows; the RT-PCR reaction was run on the thermal cycler in three stages; 2 minutes
161 at 50°C, 10 minutes at 90°C and 40 cycles of 15 seconds at 90°C followed by 1 minute at 60°C.

162 ExpressionSuite Software v1.0.3 (Life Technologies) was used to generate automated threshold values for
163 signal amplification for a majority of samples. Rarely were automated thresholds adjusted manually.
164 Amplification plots resulting in Ct values >35 were discarded, as were those plots that generated a Ct value
165 but did not display a trend of logarithmic amplification. All Ct values were exported from the
166 ExpressionSuite software and relative quantification calculations were performed in Microsoft Excel 2010.

167
168 **Digital Image analysis:** IHC images of P-cadherin, PD-L1, and 5T4 were analyzed at 20x magnification using
169 multiple software packages. The detailed implementation in each software package is described below.
170 Briefly, the traditional cell-based H-score was implemented in HALO (Version 2.3) and QuPath (Version
171 0.2.0-m2) and was calculated based on the cell-membrane localized biomarker signal. The ATM score was
172 implemented in Visiopharm (Version 2017.7.3.4069) and the pix H-score was implemented in Visiopharm
173 and HALO.

174 **HALO implementation of H-score (H-score (HALO)):** The membrane algorithm (v1.4) in HALO was used to
175 detect cells and calculate the H-score. The algorithm first deconvolves the IHC image into hematoxylin
176 and DAB channels, then detects individual cells and their subcellular compartments, i.e. nucleus and cell
177 membrane, in the image, and scores the cells as high, medium, and low based on the average DAB signal
178 associated with the cell membrane. The thresholds for high, medium, and low were determined
179 separately for each biomarker by examining the membrane-associated DAB signal across multiple images
180 pertaining to that biomarker. A separate algorithm was implemented for each biomarker in order to

181 optimize the detection and segmentation of the nucleus and cell membrane specific to that biomarker.
182 The App outputs the number of negative, high, medium and low cells, which is then used to calculate the
183 H-score that is given by

$$\text{H-score} = 100x \frac{3*(\# \text{ of high cells}) + 2*(\# \text{ of medium cells}) + \# \text{ of low cells}}{\text{Total number of cells}} \quad 1$$

184

185 **QuPath implementation of H-score (H-score (QuPath)):** QuPath is an open-source software for whole-
186 slide image analysis of histopathology data [16]. A script was written in the Groovy programming language
187 to detect cells and score them as high, medium, and low based on the average DAB signal in the cell
188 membrane. The script first deconvolves the IHC image into hematoxylin and DAB channels. A watershed-
189 based cell and membrane detection algorithm (Analyze -> Cell Analysis -> Cell + membrane detection) was
190 used to detect individual cells and identify their subcellular compartments, i.e. nucleus and cell
191 membrane. The cell detection algorithm includes a pre-processing step that involves a local background
192 subtraction by using the minimum filter. The optional median filtering step was not used. Cells that were
193 devoid of a nucleus (due to weak or missing hematoxylin staining) were excluded and the remaining cells
194 were scored as high, medium, and low based on the mean DAB signal associated with the membrane
195 compartment. The thresholds for high, medium, and low were determined separately for each biomarker.
196 A separate script was implemented for each biomarker in order to optimize the detection and
197 segmentation of the nucleus and cell membrane specific to that biomarker. The script outputs the total
198 number detected cells along with the number of high, medium, and low cells, which is then used to
199 calculate the H-score that is given in Eq. 1.

200 **ATM score:** The motivation behind the ATM score is discussed elsewhere [13]. Briefly, the idea is to use
201 all the intensity values in the DAB channel so that the final metric is independent of the choice of the
202 thresholds. Further, the ATM score is a pixel-based metric that does not depend on the detection of

203 individual cells and/or its subcellular components. Assuming 8-bit resolution for the color-deconvolved
204 biomarker channel, the ATM score is given by [13]

$$\text{ATM score} = \frac{1}{256} \sum_{k=0}^{255} PS(k), \quad 2$$

205
206 where $PS(k)$ denotes the proportion of pixels with intensity greater than or equal to k , where k takes
207 values from 0 to 255 (i.e. 2^8 grey levels). If n denotes the total number of pixels in the biomarker channel,
208 b_i denotes the biomarker intensity at the i^{th} pixel for $i = 1, \dots, n$, and $I(b_i > k)$ denotes an indicator function,
209 i.e. $I(b_i > k) = 1$ if $b_i > k$ and 0 otherwise, then the term $PS(k)$ can be written as

$$PS(k) = \frac{1}{n} (\# \text{ of pixels with values greater than } k) = \frac{1}{n} \sum_{i=1}^n I(b_i > k).$$

210
211 Substituting the above equation in Eq. 2, we have

$$\text{ATM score} = \frac{1}{256} \sum_{k=0}^{255} n^{-1} \sum_{i=1}^n I(b_i > k) = \frac{n^{-1}}{256} \sum_{i=1}^n \sum_{k=0}^{255} I(b_i > k)$$

$$= \frac{n^{-1}}{256} \sum_{i=1}^n (I(b_i > 0) + I(b_i > 1) + \dots + I(b_i > b_i - 1) + I(b_i > b_i) + \dots + I(b_i > 255))$$

$$= \frac{n^{-1}}{256} \sum_{i=1}^n (1 + 1 + \dots + 1 + 0 + \dots + 0)$$

$$= \frac{n^{-1}}{256} \sum_{i=1}^n b_i = \frac{1}{256} (\text{average value of all the pixels in the DAB channel}).$$

216 From the above equation, we see that the ATM score is a weighted average of all the pixels in the DAB
217 channel. The ATM score was implemented in Visiopharm software. The IHC image was color deconvolved
218 into hematoxylin and DAB channels. Therefore, the ATM score was calculated by taking the average
219 intensity of all DAB positive pixels and then dividing this by 256.

220 **Visiopharm implementation of pix H-score (pix H-score (VIS)):** A threshold-based detection App was used
221 to implement the pix H-score in Visiopharm. The App first deconvolves the IHC image into hematoxylin
222 and DAB channels. The App then detects and classifies DAB positive pixels as high, medium, and low, and
223 then detects the hematoxylin positive pixels. The thresholds for DAB and hematoxylin were separately
224 selected for each biomarker. The App then outputs the total area of the DAB high, DAB medium, and DAB
225 low pixels and the hematoxylin positive pixels. These values are then used to calculate the pix H-score
226 which is given by:

$$\begin{aligned} \text{pix H-score} = & \frac{100}{(\text{Total area of all DAB pixels}) + (\text{total area of hematoxylin pixels})} \\ & \times (3 * (\text{area of high DAB pixels}) + 2 * (\text{area of medium DAB pixels}) \\ & + (\text{area of low DAB pixels})). \end{aligned} \quad 3$$

229
230 In Visiopharm, the intensity-based thresholding algorithm depends on the order in which the different
231 color-deconvolved channels are specified. For instance, if a pixel contains both hematoxylin and DAB
232 signal that are above their respective threshold values for positivity and the DAB channel is first analyzed
233 followed by the hematoxylin channel, then that pixel will be labeled as positive only for the DAB channel.
234 In other words, if a pixel is found to be positive for one of the color-deconvolved channels then it is
235 excluded from any subsequent classification for the other color-deconvolved channels.

236 **HALO implementation of pix H-score (pix H-score (HALO)):** The area quantification algorithm (v2.1.3) in
237 HALO was used to calculate the pix H-score with the number of phenotypes set to 1. The algorithm
238 deconvolves the IHC image into hematoxylin and DAB channels and can detect and classify hematoxylin
239 and DAB positive pixels as high, medium, and low based on a user defined threshold. For the calculation
240 of pix H-score, a single threshold was used to detect all hematoxylin positive pixels and three separate
241 thresholds were used to detect and classify the DAB positive pixels. In HALO, these thresholds take values

242 between 0 and 1. In order to keep the thresholds implemented in Visiopharm and HALO identical, the
243 threshold values used in Visiopharm, which take values between 0 – 255, were rescaled to take values
244 between 0 and 1 and these were then used in HALO. Unlike Visiopharm, HALO keeps track of the detected
245 pixels in the DAB and hematoxylin channels separately. Consequently, pixels that contain both DAB and
246 hematoxylin signal that are above the thresholds will be accounted for in both the hematoxylin and DAB
247 channels. In order to mimic the Visiopharm implementation of pix H-score, we define a third channel,
248 which is denoted as phenotype 1 channel in HALO that pertains to pixels that are positive for hematoxylin
249 but negative for DAB. This phenotype 1 channel will contain pixels that are analogous to the hematoxylin
250 positive pixels detected in the Visiopharm implementation of pix H-score algorithm. The algorithm
251 outputs the area high, medium, and low pixels in the DAB channel, the area of positive pixels in the
252 phenotype 1 channels, which is used to as an estimate of the total area of pixels containing only the
253 hematoxylin signal. These values are then used in Eq. 3 to calculate the pix H-score.

254 **Statistical analysis:** Spearman's rank correlation coefficient was calculated to assess the correlation
255 between the DIA endpoint and biomarker abundance. The William's t test was used to test for significant
256 difference between a pair of dependent correlation coefficients [17, 18].

257 **Spatial resampling analysis:** For each biomarker, an empirical resampling procedure was performed on
258 every whole-slide IHC image. The viable tissue region was sampled by non-overlapping circular regions of
259 radius 0.8 mm (Figure 6A). For each region, the area of DAB high, DAB medium, DAB low, and hematoxylin
260 positive pixels were determined using Visiopharm. The results were exported to MATLAB (Mathworks,
261 Natick, MA) for subsequent analysis. For every IHC image, N different circular regions were randomly
262 selected ($N = 1 - 50$), and a regional pix H-score was calculated using the area of DAB high pixels, DAB
263 medium pixels, DAB low pixels, and hematoxylin positive pixels that were summed from the N circular
264 regions. This procedure is repeated N_{iter} times with replacement ($N_{iter} = 100$ for all the biomarkers). Then

265 for each iteration $k = 1, \dots, N_{iter}$, the Spearman correlation coefficient $C(N, k)$ is computed between the
266 regional pix H-score and the corresponding pathologist H-score. The average Spearman correlation
267 coefficient for each value of N is computed using the formula

$$C_{av}(N) = \frac{1}{N_{iter}} \sum_{k=1}^{N_{iter}} C(N, k).$$

268

269

270 **RESULTS**

271 **DIA algorithms for P-cadherin quantification**

272 IHC images for P-cadherin (Figure 2A) showed strong immunoreactivity at the cell membrane and in the
273 cytoplasm, which was consistent with prior reports [19, 20]. Spearman's correlation analysis of the
274 membrane H-scores of the 30 cases immunolabeled for P-cadherin, as assessed by a board-certified
275 pathologist (see Supplementary Table 1), and NanoString nCounter values for P-cadherin mRNA transcript
276 from serial sections of the same cases had a correlation coefficient of 0.81, $p < 0.0001$ (Figure 2B).

277 When compared to the P-cadherin pathologist H-score, all P-cadherin DIA endpoints yielded positive
278 correlations (Figures 3A-3E). The correlation with the ATM score (Figure 3C) and pix H-score (Figures 3D
279 and 3E) were higher than the correlations with the DIA based H-scores (Figures 3A and 3B). More
280 specifically, the Spearman's correlation coefficient for HALO and QuPath DIA H-scores were 0.5 ($p = 0.005$)
281 and 0.39 ($p = 0.03$), respectively, whereas the Spearman's correlation coefficient for the ATM score, the
282 VIS pix H-score and the HALO pix H-score were 0.78 ($p < 0.001$), 0.77 ($p < 0.0001$) and 0.88 ($p < 0.0001$),
283 respectively. When compared to the P-cadherin transcript, all DIA endpoints similarly yielded positive
284 correlations (Figures 3F-3J), with the pix H-score exhibiting the highest Spearman's correlation coefficient
285 (Figures 3I and 3J; $\rho = 0.83$ and $\rho = 0.81$, respectively, for VIS and HALO pix H-score; $p < 0.0001$) followed
286 by the ATM score (Figure 3H; $\rho = 0.62$, $p < 0.0001$) and the DIA H-scores (Figures 3F and 3G; $\rho = 0.5$, $p =$
287 0.005 for HALO and $\rho = 0.45$, $p = 0.01$ for QuPath).

288 We next investigated whether the differences in the Spearman correlation coefficients for the various DIA
289 endpoints are statistically significant. Table 1 shows the results of our statistical analysis where we carried
290 out pairwise comparisons of the correlation coefficients for different DIA endpoints obtained from P-
291 cadherin IHC images. Our analysis shows that the correlation coefficient between the pix H-score and
292 either of the biomarker abundance endpoints (pathologist H-score and P-cadherin transcript) is

293 significantly higher than the correlation coefficient between DIA based H-scores and biomarker
294 abundance endpoints. This suggests that for the P-cadherin dataset, the pix H-score is a better DIA metric
295 to quantify biomarker abundance over traditional DIA based H-score. In the case of the ATM score, we
296 observe a mixed result in that the correlation coefficient between pix H-score and P-cadherin transcript
297 is significantly higher than the correlation coefficient between ATM score and P-cadherin transcript,
298 whereas statistical significance is lost when we consider the pathologist H-score as the reference for
299 biomarker abundance (Table 1). We also compared the two DIA based H-scores. We found no significant
300 difference in the Spearman's correlation coefficient between QuPath H-score and biomarker abundance
301 endpoints versus HALO H-score and biomarker abundance endpoints (Table 1). Similarly, we found no
302 significant difference in the correlation coefficients for the HALO and VIS implementations of the pix H-
303 score for P-cadherin.

304 **DIA algorithms for PD-L1 quantification**

305 IHC images for PD-L1 (Figure 2C) showed strong immunoreactivity at the cell membrane and minimal to
306 no cytoplasmic staining, which was consistent with prior reports [19, 20]. Spearman's correlation analysis
307 of the membrane H-scores of the 24 cases immunolabeled for PD-L1, as assessed by a board-certified
308 pathologist (see Supplementary Table 1), and NanoString nCounter values for PD-L1 mRNA transcript
309 from serial sections of the same cases had a correlation coefficient of 0.91, $p < 0.0001$ (Figure 2D).

310 When compared to the pathologist H-score, all DIA endpoints yielded positive correlations (Figures 4A-
311 4E). The Spearman's correlation coefficient for the HALO H-score, QuPath H-score, ATM score, VIS pix H-
312 score and HALO pix H-score with respect to the pathologist H-score were 0.69 ($p = 0.0002$), 0.74 ($p < 0.0001$),
313 0.55 ($p = 0.005$), 0.76 ($p < 0.0001$) and 0.71 ($p < 0.0001$), respectively. When compared to the PD-L1
314 transcript, all DIA endpoints similarly yielded positive correlations (Figures 4F-4J). The Spearman's
315 correlation coefficient for the HALO H-score, QuPath H-score, ATM score, VIS pix H-score and HALO pix H-

316 score with respect to PD-L1 transcript were 0.73 ($p < 0.0001$), 0.75 ($p < 0.0001$), 0.55 ($p = 0.005$), 0.79
317 ($p < 0.0001$) and 0.79 ($p < 0.0001$), respectively.

318 Statistical analysis of the Spearman's correlation coefficients revealed that there is no significant
319 difference in the correlation coefficient between DIA based H-scores and PD-L1 biomarker abundance
320 endpoints versus the correlation coefficient between pix H-score and PD-L1 biomarker abundance
321 endpoints (Table 2). This shows that the performance of pix H-score is analogous to that of the DIA based
322 H-score which is in contrast with our observations for P-cadherin. Also, there was no significant difference
323 in Spearman's correlation coefficient between HALO and QuPath implementations of the H-score, which
324 is analogous to what we observed for P-cadherin. In addition, we observed that there was no significant
325 difference between the HALO and Visiopharm implementations of the pix H-score for PD-L1. Spearman's
326 correlation coefficients between the pix H-score and PD-L1 biomarker abundance endpoints were mostly
327 significantly higher than Spearman's correlation coefficients between ATM score and PD-L1 biomarker
328 abundance endpoints (Table 2). Although both the pix H-score and the ATM score are pixel-based
329 algorithms, the higher Spearman's correlation coefficient for the pix H-score suggests that this algorithm
330 is superior to the ATM score in estimating biomarker abundance for PD-L1.

331 **DIA algorithms for 5T4 quantification**

332 IHC images for 5T4 (Figure 2E) showed strong immunoreactivity at the cell membrane with limited
333 cytoplasmic staining, which was consistent with prior reports [15]. Spearman's correlation of the
334 membrane H-scores of the 21 cases immunolabeled for 5T4, as assessed by a board-certified pathologist
335 (see Supplementary Table 1), and qRT-PCR values for 5T4 mRNA transcript from serial sections of the same
336 cases had a ρ value of 0.61, $p = 0.003$ (Figure 2F).

337 When compared to the pathologist H-score, all DIA endpoints yielded positive correlations (Figures 5A-
338 5E). The Spearman's correlation coefficient for the HALO H-score, QuPath H-score, ATM score, VIS pix H-

339 score and HALO pix H-score with respect to the pathologist H-score were 0.75 ($p < 0.0001$), 0.79 ($p < 0.0001$),
340 0.76 ($p < 0.0001$), 0.83 ($p < 0.0001$) and 0.82 ($p < 0.0001$), respectively. When compared to the 5T4
341 transcript, all DIA endpoints similarly yielded positive correlations (Figures 5F-5J). The Spearman's
342 correlation coefficient for the HALO H-score, Qupath H-score, ATM score, VIS pix H-score and HALO pix H-
343 score with respect to 5T4 transcript were 0.74 ($p < 0.0001$), 0.55 ($p = 0.01$), 0.69 ($p = 0.0007$), 0.76
344 ($p < 0.0001$) and 0.74 ($p = 0.0001$), respectively.

345 Statistical analysis of the Spearman's correlation coefficients revealed that there is no significant
346 difference in the correlation coefficient between each of the DIA based endpoints and pathologist H-score
347 (Table 3). An analogous behavior was also observed for the correlation coefficient between each of the
348 DIA based endpoints and 5T4 transcript except for the QuPath H-score. Specifically, the correlation
349 between QuPath H-score and 5T4 transcript was significantly lower than the correlation between the
350 HALO H-score or the pix H-score endpoints and 5T4 transcript (Table 3). Finally, we note that there is no
351 significant difference in the correlation coefficient between the HALO and Visiopharm implementations
352 of the pix H-score and either of the biomarker abundance endpoints for 5T4. These results suggest that
353 the pix H-score algorithm has comparable performance to the other DIA algorithms to quantify biomarker
354 abundance for 5T4.

355 **Effect of spatial sampling on pix H-score**

356 We next investigated the robustness of the pix H-score when it is calculated from select regions within
357 the tissue section as opposed to the entire tumor resection. For this purpose, a statistical sampling
358 procedure known as bootstrapping needs to be performed. However, technical limitations in Visiopharm
359 and HALO software packages precluded us from implementing a formal bootstrapping procedure.
360 Therefore, we resorted to an empirical resampling approach (see Methods for details) wherein for a given
361 biomarker each tumor resection was divided into non-overlapping circular regions (Figure 6A). N different

362 circular regions (N ranging from 1 to 50) were randomly selected, and a regional pix H-score was computed
363 from these circular regions. Then the Spearman's correlation coefficient between the pathologist H-score
364 and the regional pix H-score was computed for that biomarker. This procedure was repeated 100 times
365 for all the tumor resections pertaining to that biomarker, and the average Spearman correlation
366 coefficient from 100 iterations was then plotted as a function of the number of circular regions N.

367 Figures 6B, 6C and 6D show the behavior of the average Spearman's correlation coefficient for PD-L1, P-
368 cadherin and 5T4, respectively, between pathologist H-score and the regional pix H-score as a function of
369 the number of circular regions from which the regional pix H-score was calculated. For all the biomarkers,
370 we see that for fewer than five circular regions the average Spearman correlation coefficient between the
371 regional pix H-score and pathologist H-score is consistently smaller than the Spearman's correlation
372 coefficient between the whole-slide pix H-score and pathologist H-score (shown by the red dashed line).
373 When 10 or more circular regions are sampled the average Spearman's correlation coefficient for the
374 regional pix H-score starts to plateau out and reaches a steady state. In the case of PD-L1, the plateau
375 region converges with the Spearman's correlation coefficient between the whole-slide pix H-score and
376 pathologist H-score (Figure 6B). In contrast, for P-cadherin 5T4 the plateau region is slightly lower than
377 the Spearman's correlation coefficient for the whole-slide pix H-score (Figures 6C and 6D). A similar
378 behavior is also observed when biomarker mRNA levels are used as the reference ground truth data in
379 the Spearman's correlation coefficient calculation (data not shown).

380

381

382

383 **DISCUSSION**

384 Robust quantification of biomarker expression in tissue sections is a critical need in many diagnostic and
385 investigative pathology workflows. Our motivation to develop a new digital image analysis metric was
386 driven by the need to automate the process of manual scoring by a pathologist. Digital image analysis
387 holds the promise to offer a fast, objective, and reproducible strategy to quantify biomarker expression
388 from histopathology images. In this manuscript, we introduced an unsupervised algorithm, the pix H-
389 score. With it we quantified P-cadherin, PD-L1, and 5T4 signals in immunolabeled FFPE sections of human
390 tumors and found good correlation between the digitally-analyzed IHC signals and manual (visual) signal
391 quantitation as performed by a board certified pathologist. As pathologist scoring is known to be
392 susceptible to intra- and inter-observer variability, we also used biomarker mRNA level as an orthogonal
393 measurement of biomarker abundance to validate the pix H-score. Our observation that there was good
394 concordance between both digital and visual IHC signal quantitation and mRNA transcript abundance for
395 each analyte not only demonstrated the robust nature of the pix H-score algorithm but also validated the
396 pathologist scores.

397 There are two basic approaches to quantifying biomarker expression from histology images. One
398 approach utilizes cell segmentation and quantifies markers per unit cell whereas a second approach
399 avoids cell segmentation and quantifies markers per unit pixel. In this manuscript, we compared both
400 approaches to quantify biomarker levels from immunohistochemistry images. Unlike the H-score and the
401 Allred score, the pix H-score is a pixel-based algorithm that does not rely on the identification of individual
402 cells and their subcellular compartments. This reduces the computational complexity of the pix H-score
403 and renders its implementation in two different software packages as relatively straightforward.

404 In our case, the IHC assay for each biomarker was carried out using a different brand of instrument (PD-
405 L1 – Ventana, P-cadherin – DAKO, and 5T4 – Leica Bond RX). Similarly, the slides were scanned using

406 different whole-slide scanners (PD-L1 and P-cadherin - different Aperio AT2 scanners, and 5T4 –
407 Hamamatsu NanoZoomer). These differences could introduce variations in the colorimetric composition
408 of the IHC images that can impact downstream image analysis. Our observation that the Visiopharm and
409 the HALO versions of pix H-score exhibited similar performance suggests that the pix H-score is a robust
410 algorithm for estimating IHC biomarker abundance in whole-slide images. This is especially relevant due
411 to the proprietary nature of these software packages which precludes users from understanding several
412 technical aspects of the image analysis workflow. For instance, the specific details regarding the color
413 deconvolution algorithm, which is a key pre-processing step, are not accessible to the user in either
414 Visiopharm or HALO. Consequently, while implementing the pix H-score we did not know how similar the
415 output of the color deconvolution step (i.e. hematoxylin and DAB channels) would be in the two software
416 packages.

417 An important question arises as to why the DIA based H-score exhibited very different performance for P-
418 cadherin but not for PD-L1. The H-score algorithm relied on the detection of individual cells and their
419 subcellular compartments to quantify biomarker levels. Although this task may seem relatively
420 straightforward for a human observer, nucleus/cell-membrane detection and segmentation are
421 challenging image processing problems especially when applied to whole-slide image analysis where there
422 can be considerable variability in the intensity and the sub-cellular localization pattern of the biomarker
423 of interest [21, 22]. In our case, the latter could be a contributing factor since in the P-cadherin IHC images
424 the biomarker signal was localized to both the cell membrane and cytoplasm whereas in the PD-L1 IHC
425 images the biomarker signal was predominantly localized to the cell membrane. Consequently, this may
426 partly explain the reason why for P-cadherin the performance of the DIA H-score was consistently lower
427 than that of the pix H-score whereas for PD-L1 the performance of the DIA H-score was comparable to
428 that of the pix H-score. Not surprisingly others have also reported similar challenges in automated analysis
429 of membrane-localized biomarker signal [23]. This may also partly explain our observation for 5T4 where

430 the correlation between QuPath H-score and 5T4 transcript was lower than the correlation between pix
431 H-score and 5T4 transcript. More specifically, while 5T4 immunoreactivity is predominantly membranous,
432 there is still detectable cytoplasmic signal in the tumor cells which can affect the quantification of the DIA
433 based H-score.

434 A similar question also arises for the ATM score which, unlike the H-score, is a pixel-based algorithm but
435 also exhibited very different performance for P-cadherin but not for PD-L1 and 5T4. By definition, the ATM
436 score is proportional to the average intensity of the biomarker in the DAB channel. This is calculated by
437 taking all pixels in the DAB channel including pixels that are negative for the biomarker. When the
438 averaging is performed on a whole-slide image, this can significantly dilute the contribution from pixels
439 that are positive for the biomarker resulting in poor performance in predicting biomarker abundance from
440 the IHC image. In contrast, the pix H-score only considers pixels with a valid biomarker signal as DAB
441 positive pixels (based on a user defined threshold). As a result, the pix H-score can robustly estimate
442 biomarker abundance the IHC image. This difference also explains in part the reason for the limited range
443 of values taken by the ATM score when compared to the pix H-score. Specifically, the ATM score for P-
444 cadherin, PD-L1, and 5T4 took values in the range of 24 to 77, 8 to 33, and 11 to 49, respectively. In
445 contrast the pix H-score for P-cadherin, PD-L1, and 5T4 took values in the range of 20 to 207, 1 to 131,
446 and 3 to 170, respectively. The latter values are more comparable to the pathologist H-score, which for P-
447 cadherin, PD-L1, and 5T4 ranged from 17 to 298, 0 to 225, and 0 to 224, respectively.

448 The application of deep learning methodology for nucleus and cell membrane segmentation holds
449 significant promise as it has been shown to have improved performance over traditional algorithms [24].
450 However, deep learning methods are supervised approaches that require a substantial amount of training
451 data and extensive validation. In many practical applications, generating such large training datasets is
452 not feasible and algorithm validation can be time consuming. In this regard, the pix H-score algorithm

453 introduced here provides a simple yet robust strategy to quantify biomarker expression even from small
454 datasets, as demonstrated here, and can be implemented within a very short timeframe. An interesting
455 follow up study would be to compare the performance of the pix H-score algorithm with deep learning
456 based, scoring approaches.

457 We note that while our results are encouraging and show the potential for the pix H-score in scoring
458 membrane biomarkers, the algorithm can benefit from additional validation for other biomarkers. Also,
459 the effect of pre-analytical variables (e.g., cold ischemia time, age of unstained cut slides, etc.) on the
460 performance of the pix H-score needs to be investigated. In addition, the effect of stain variation needs
461 to be explored, which is known to be a notable source of variability in histopathology data. In our current
462 work stain normalization was not necessary, likely due to the small batch size of our datasets which did
463 not exhibit significant colorimetric variability. Although not shown here, we expect the pix H-score to also
464 be applicable to immunofluorescence images. In conclusion, we anticipate the pix H-score to be a useful
465 addition to the digital image analysis toolbox for a fast, reproducible and objective strategy to quantify
466 biomarker expression from immunolabeled tissue sections.

467

468 **ACKNOWLEDGEMENTS:** We thank Shawn O'Neil and Timothy Affolter for critical reading of the
469 manuscript.

470

471

472

473

474

475

476 **REFERENCES**

- 477 1. Aeffner, F., et al., *Introduction to Digital Image Analysis in Whole-slide Imaging: A White Paper*
478 *from the Digital Pathology Association*. J Pathol Inform, 2019. **10**: p. 9.
- 479 2. Meyerholz, D.K. and A.P. Beck, *Principles and approaches for reproducible scoring of tissue stains*
480 *in research*. Lab Invest, 2018. **98**(7): p. 844-855.
- 481 3. Aeffner, F., et al., *Commentary: Roles for Pathologists in a High-throughput Image Analysis*
482 *Team*. Toxicol Pathol, 2016. **44**(6): p. 825-34.
- 483 4. Brunstrom, H., et al., *PD-L1 immunohistochemistry in clinical diagnostics of lung cancer: inter-*
484 *pathologist variability is higher than assay variability*. Mod Pathol, 2017. **30**(10): p. 1411-1421.
- 485 5. Gomes, D.S., et al., *Inter-observer variability between general pathologists and a specialist in*
486 *breast pathology in the diagnosis of lobular neoplasia, columnar cell lesions, atypical ductal*
487 *hyperplasia and ductal carcinoma in situ of the breast*. Diagn Pathol, 2014. **9**: p. 121.
- 488 6. Hirsch, F.R., et al., *PD-L1 Immunohistochemistry Assays for Lung Cancer: Results from Phase 1 of*
489 *the Blueprint PD-L1 IHC Assay Comparison Project*. J Thorac Oncol, 2017. **12**(2): p. 208-222.
- 490 7. Rimm, D.L., et al., *A Prospective, Multi-institutional, Pathologist-Based Assessment of 4*
491 *Immunohistochemistry Assays for PD-L1 Expression in Non-Small Cell Lung Cancer*. JAMA Oncol,
492 2017. **3**(8): p. 1051-1058.
- 493 8. Rizzardi, A.E., et al., *Quantitative comparison and reproducibility of pathologist scoring and*
494 *digital image analysis of estrogen receptor beta2 immunohistochemistry in prostate cancer*.
495 Diagn Pathol, 2016. **11**(1): p. 63.

- 496 9. Barnes, M., et al., *Whole tumor section quantitative image analysis maximizes between-*
497 *pathologists' reproducibility for clinical immunohistochemistry-based biomarkers.* Lab Invest,
498 2017. **97**(12): p. 1508-1515.
- 499 10. Tsao, M.S., et al., *PD-L1 Immunohistochemistry Comparability Study in Real-Life Clinical Samples:*
500 *Results of Blueprint Phase 2 Project.* J Thorac Oncol, 2018. **13**(9): p. 1302-1311.
- 501 11. Stalhammar, G., et al., *Digital image analysis outperforms manual biomarker assessment in*
502 *breast cancer.* Mod Pathol, 2016. **29**(4): p. 318-29.
- 503 12. Aeffner, F., et al., *The Gold Standard Paradox in Digital Image Analysis: Manual Versus*
504 *Automated Scoring as Ground Truth.* Arch Pathol Lab Med, 2017. **141**(9): p. 1267-1275.
- 505 13. Choudhury, K.R., et al., *A robust automated measure of average antibody staining in*
506 *immunohistochemistry images.* J Histochem Cytochem, 2010. **58**(2): p. 95-107.
- 507 14. Camp, R.L., G.G. Chung, and D.L. Rimm, *Automated subcellular localization and quantification of*
508 *protein expression in tissue microarrays.* Nat Med, 2002. **8**(11): p. 1323-7.
- 509 15. Pirie-Shepherd, S.R., et al., *Detecting expression of 5T4 in CTCs and tumor samples from NSCLC*
510 *patients.* PLoS One, 2017. **12**(7): p. e0179561.
- 511 16. Bankhead, P., et al., *QuPath: Open source software for digital pathology image analysis.* Sci Rep,
512 2017. **7**(1): p. 16878.
- 513 17. Steiger, J.H., *Tests for comparing elements of a correlation matrix.* Psychological Bulletin, 1980.
514 **87**(2): p. 245-251.
- 515 18. Williams, E.J., *The Comparison of Regression Variables.* Journal of the Royal Statistical Society:
516 Series B (Methodological), 1959. **21**(2): p. 396-399.
- 517 19. Kovacs, A., J. Dhillon, and R.A. Walker, *Expression of P-cadherin, but not E-cadherin or N-*
518 *cadherin, relates to pathological and functional differentiation of breast carcinomas.* Mol Pathol,
519 2003. **56**(6): p. 318-22.

- 520 20. Paredes, J., et al., *P-cadherin overexpression is an indicator of clinical outcome in invasive breast*
521 *carcinomas and is associated with CDH3 promoter hypomethylation*. Clin Cancer Res, 2005.
522 **11**(16): p. 5869-77.
- 523 21. Irshad, H., et al., *Methods for nuclei detection, segmentation, and classification in digital*
524 *histopathology: a review-current status and future potential*. IEEE Rev Biomed Eng, 2014. **7**: p.
525 97-114.
- 526 22. Xing, F. and L. Yang, *Robust Nucleus/Cell Detection and Segmentation in Digital Pathology and*
527 *Microscopy Images: A Comprehensive Review*. IEEE Rev Biomed Eng, 2016. **9**: p. 234-63.
- 528 23. Lopes, N., et al., *Digital image analysis of multiplex fluorescence IHC in colorectal cancer*
529 *recognizes the prognostic value of CDX2 and its negative correlation with SOX2*. Lab Invest, 2020.
530 **100**(1): p. 120-134.
- 531 24. Caicedo, J.C., et al., *Evaluation of Deep Learning Strategies for Nucleus Segmentation in*
532 *Fluorescence Images*. Cytometry A, 2019. **95**(9): p. 952-965.

533

534

535 **Figure legends**

536 **Figure 1: Overview of the different scoring algorithms.** Panel A shows the traditional cell-based H-score,
537 panel B shows the average threshold method (ATM) score, and panel C shows the pix H-score.

538 **Figure 2. P-cadherin, PD-L1 and 5T4 IHC datasets.** Panels A, C and E show representative images at 20x
539 magnification with varying levels of P-cadherin, PD-L1 and 5T4 expression, respectively, in tumor
540 resections. Panels B, D and F show the plot of the pathologist H-score versus mRNA transcript level for P-
541 cadherin (n = 30 cases), PD-L1 (n = 24 cases) and 5T4 (n = 21 cases), respectively. The panels also show
542 Spearman's correlation coefficient along with the p-value and 95% confidence interval.

543 **Figure 3. Performance of DIA endpoints obtained from P-cadherin IHC images.** Panels A through E show
544 the plots of different DIA endpoints versus pathologist H-score for a cohort of 30 head and neck cancer
545 resections. Panels F through J show the plots of different DIA endpoints versus P-cadherin mRNA
546 transcript for the same 30 cases. Each panel also shows the Spearman's correlation coefficient between
547 the two quantities plotted in that panel along with the p-value and the 95% confidence interval.

548 **Figure 4. Performance of DIA endpoints obtained from PD-L1 IHC images.** Panels A-E show plots of the
549 different DIA endpoints as a function of the pathologist H-score, while panels F-J show the same as a
550 function of PD-L1 mRNA transcript for a cohort of 24 lung cancer resections. All panels show the
551 Spearman's correlation coefficient between the two quantities plotted in that panel along with the p-
552 value and the 95% confidence interval.

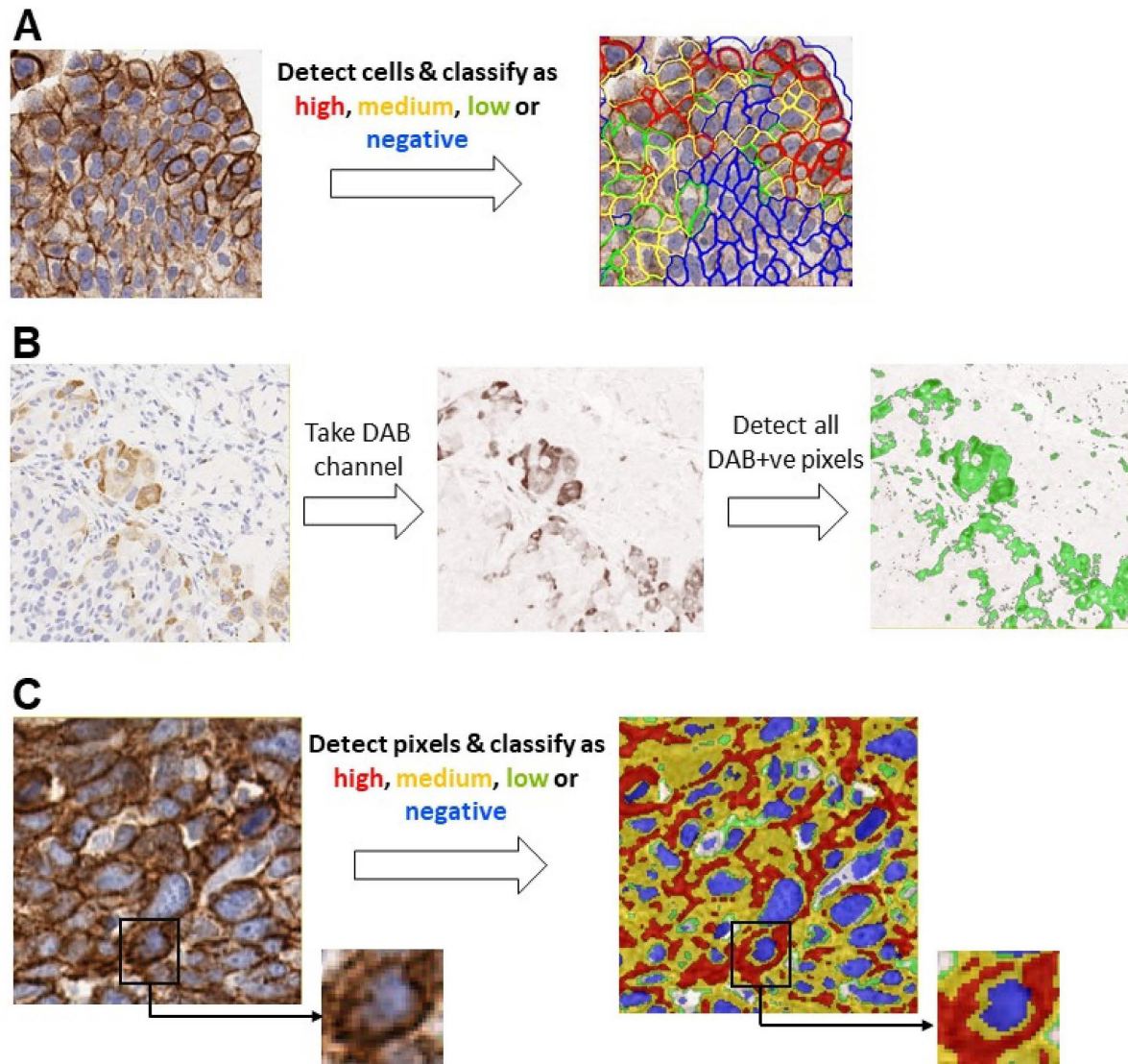
553 **Figure 5. Performance of DIA endpoints obtained from 5T4 IHC images.** Panels A-E show plots of the
554 different DIA endpoints as a function of the pathologist H-score, while panels F-J show the same as a
555 function of 5T4 mRNA transcript for a cohort of 21 lung cancer resections. All panels show the Spearman's

556 correlation coefficient between the two quantities plotted in that panel along with the p-value and the
557 95% confidence interval.

558 **Figure 6. Empirical approach to assess robustness of pix H-score to spatial sampling.** Panel A shows the
559 breakup of the tumor resection into non overlapping circular regions. Panels B, C and D show the results
560 of the bootstrap analysis for PD-L1, P-cadherin and 5T4, respectively, where the average Spearman's
561 Correlation coefficient between the regional pix H-score estimate from N circular regions and pathologist
562 H-score is plotted as a function of the number of circular regions, where N varies from 1 to 50. The red
563 dashed line shows the Spearman's correlation coefficient between whole-slide Pix H-score and
564 pathologist H-score for that biomarker. Error bars indicate \pm SEM.

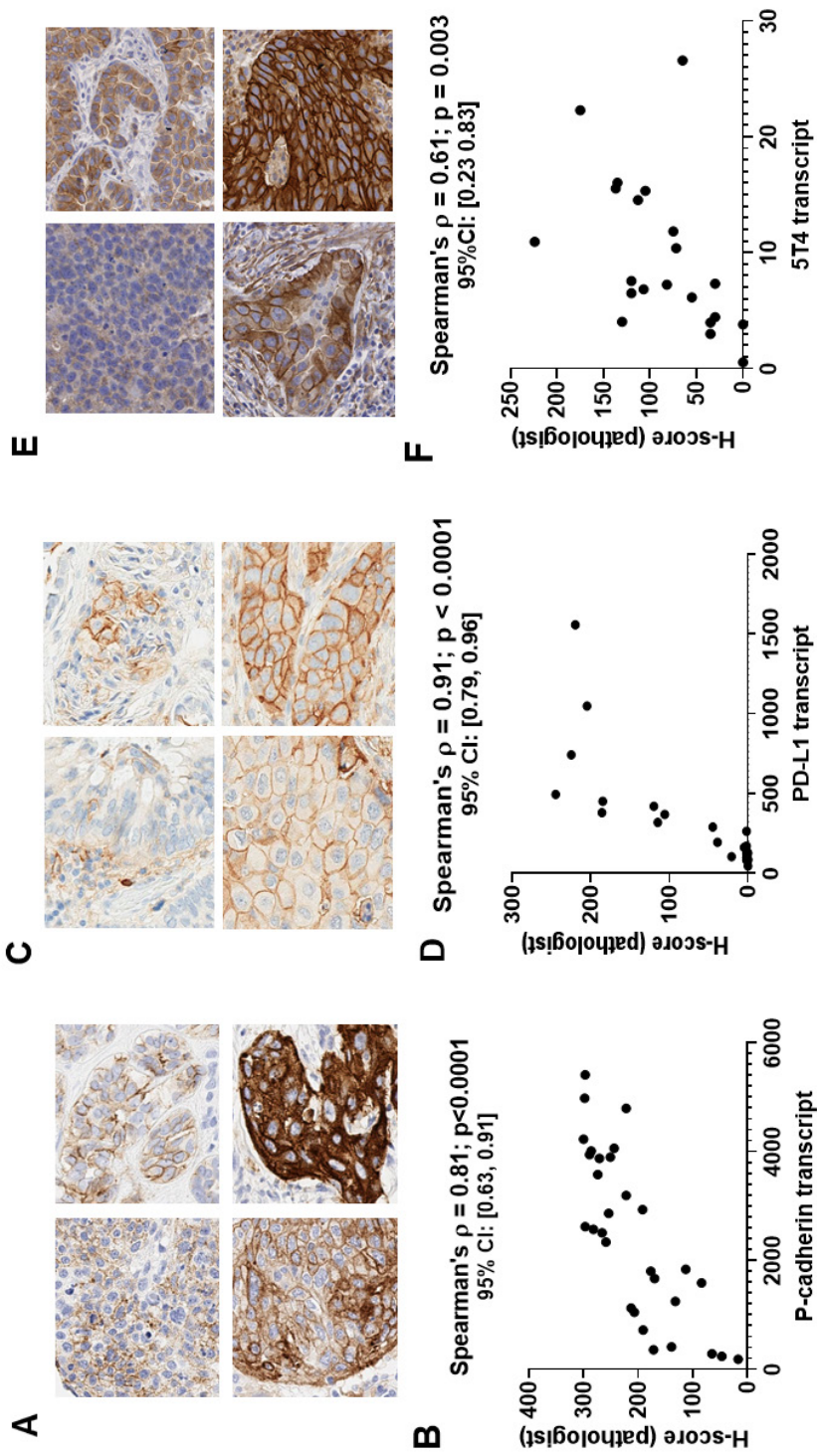
565

566 **Figure 1**

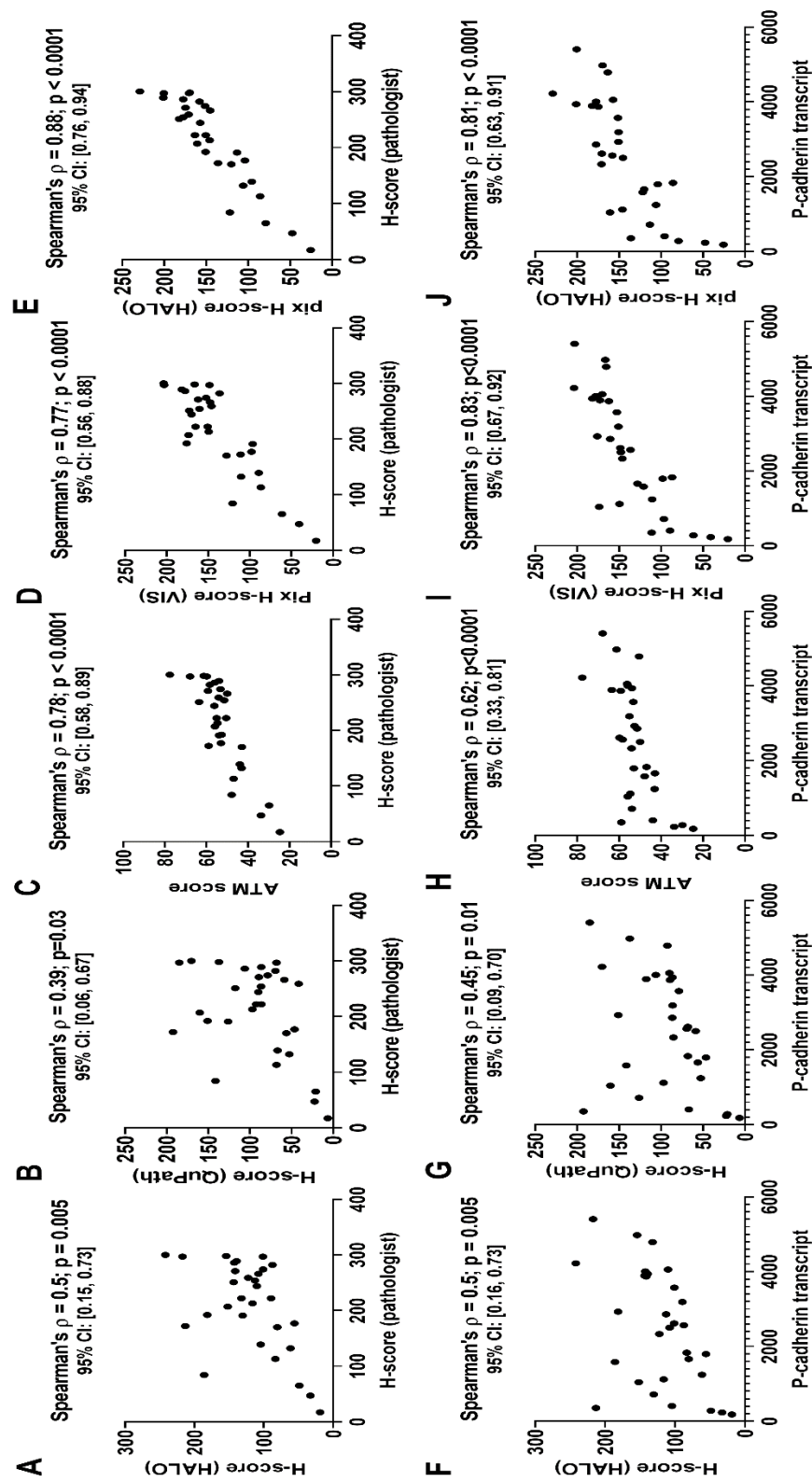


567

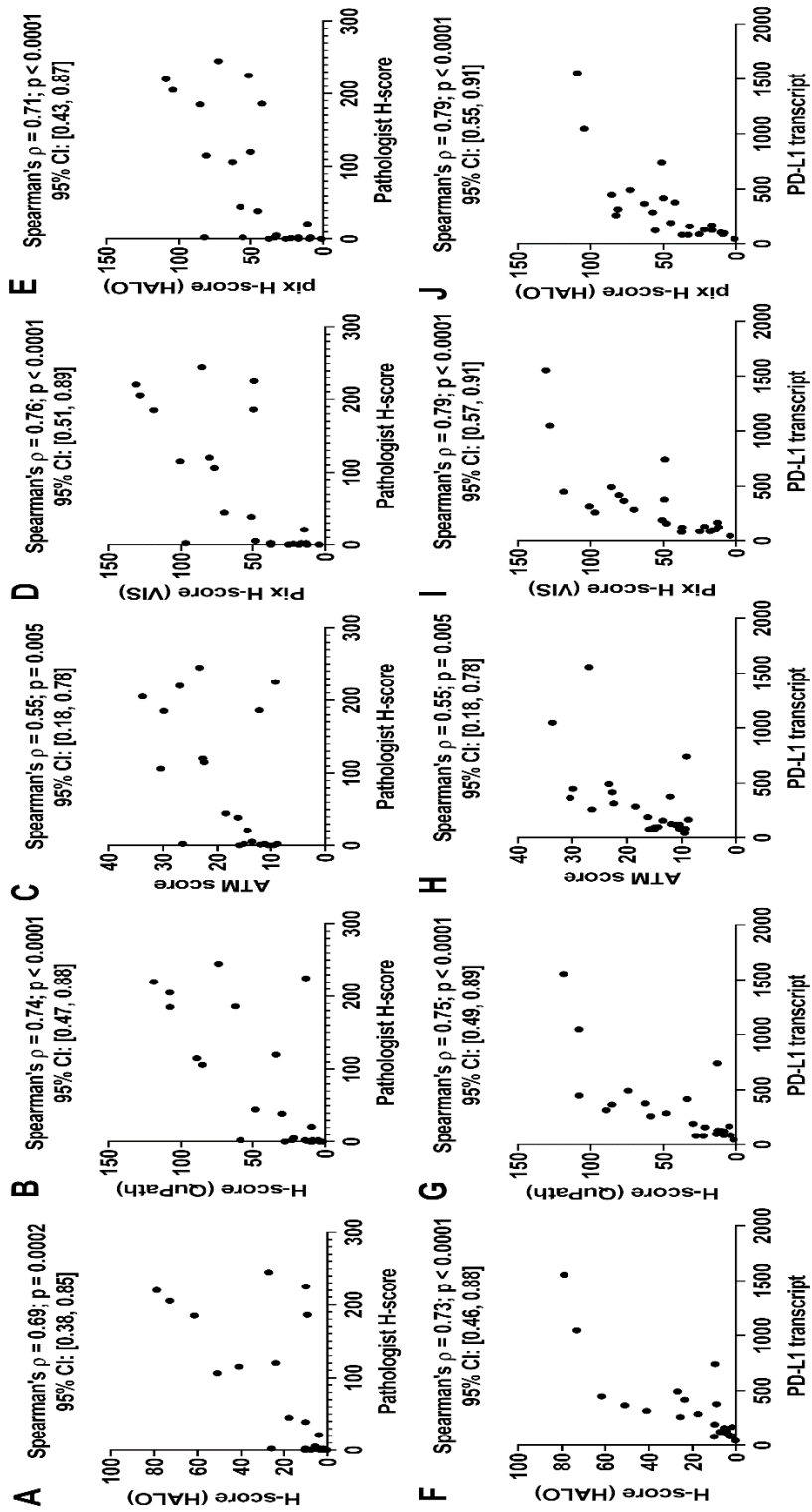
568 **Figure 2**



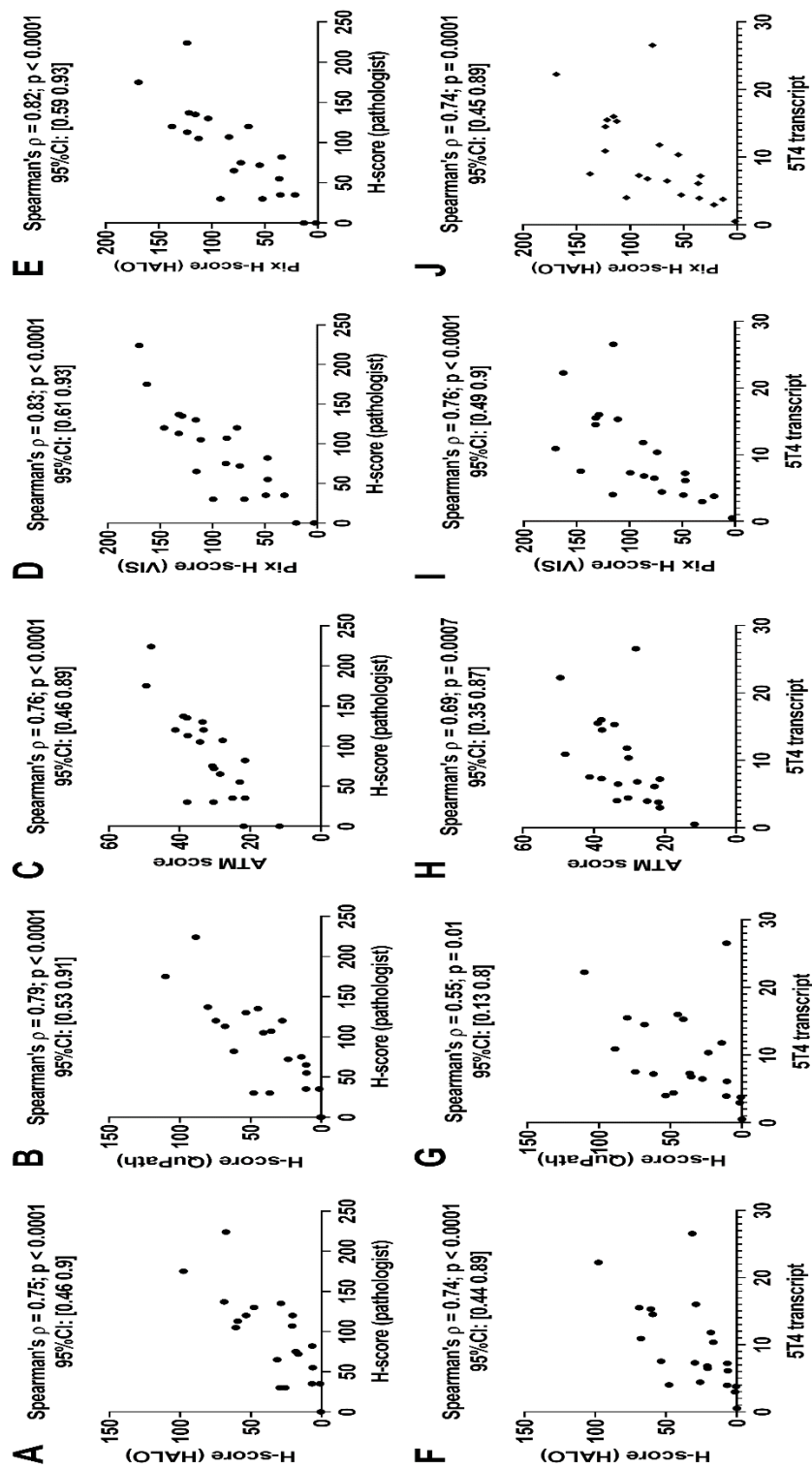
570 **Figure 3**



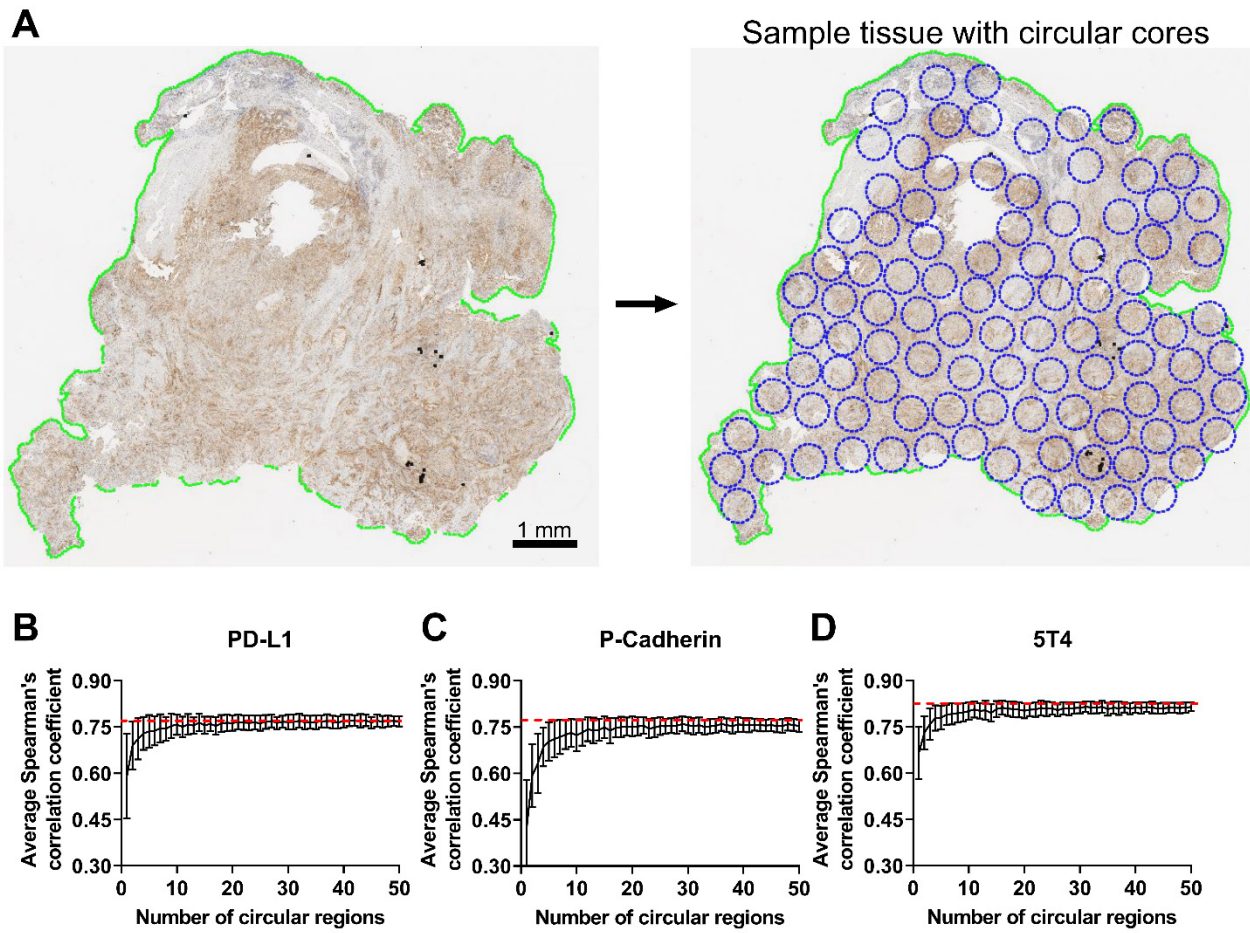
572 **Figure 4**



574 **Figure 5**



576 **Figure 6**



577

Statistical analysis of correlation coefficients for P-cadherin						
	ρ_{12}	ρ_{13}	z-score	p-value	Result	
Comparing pairwise correlations between DIA endpoint and Pathologist H-score						
$\rho(\text{Path H-score, H-score HALO})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.50	0.39	1.48	0.15	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, pix H-score HALO})$	0.77	0.88	-2.5	0.01	Significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, H-score HALO})$	0.77	0.50	3.01	0.005	Significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.77	0.39	3.87	0.0006	Significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, ATM score})$	0.77	0.78	-0.12	0.90	No significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, H-score HALO})$	0.88	0.50	5.21	1.7e-05	Significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.88	0.39	5.34	1.2E-05	Significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, ATM score})$	0.88	0.78	1.64	0.11	No significant difference	
Comparing pairwise correlations between DIA endpoint and P-cadherin transcript						
$\rho(\text{Pcad transcript, H-score HALO})$ vs $\rho(\text{Pcad transcript, H-score QuPath})$	0.50	0.45	0.79	0.43	No significant difference	
$\rho(\text{Pcad transcript, pix H-score VIS})$ vs $\rho(\text{Pcad transcript, pix H-score HALO})$	0.83	0.81	0.50	0.62	No significant difference	
$\rho(\text{Pcad transcript, pix H-score VIS})$ vs $\rho(\text{Pcad transcript, H-score HALO})$	0.83	0.5	4.19	0.0002	Significant difference	
$\rho(\text{Pcad transcript, pix H-score VIS})$ vs $\rho(\text{Pcad transcript, H-score QuPath})$	0.83	0.45	4.49	0.0001	Significant difference	
$\rho(\text{Pcad transcript, pix H-score VIS})$ vs $\rho(\text{Pcad transcript, ATM score})$	0.83	0.62	2.45	0.02	Significant difference	
$\rho(\text{Pcad transcript, pix H-score HALO})$ vs $\rho(\text{Pcad transcript, H-score HALO})$	0.81	0.5	3.31	0.002	Significant difference	
$\rho(\text{Pcad transcript, pix H-score HALO})$ vs $\rho(\text{Pcad transcript, H-score QuPath})$	0.81	0.45	3.22	0.003	Significant difference	
$\rho(\text{Pcad transcript, pix H-score HALO})$ vs $\rho(\text{Pcad transcript, ATM score})$	0.81	0.62	2.37	0.025	Significant difference	

Table 1. Table lists the results of William's t-test to test for significant difference in the Spearman correlation coefficients between P-cadherin transcript or pathologist H-score and different DIA endpoints.

Statistical analysis of correlation coefficients for PD-L1						
	ρ_{12}	ρ_{13}	z-score	p-value	Result	
Comparing pairwise correlations between DIA endpoint and Pathologist H-score						
$\rho(\text{Path H-score, H-score HALO})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.69	0.74	-1.04	0.31	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, pix H-score HALO})$	0.76	0.72	0.93	0.36	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, H-score HALO})$	0.76	0.69	1.23	0.23	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.76	0.74	0.34	0.74	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, ATM score})$	0.76	0.55	2.58	0.02	Significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, H-score HALO})$	0.72	0.69	0.44	0.67	No significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.72	0.74	-0.33	0.74	No significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, ATM score})$	0.72	0.55	1.54	0.14	No significant difference	
Comparing pairwise correlations between DIA endpoint and PDL1 transcript						
$\rho(\text{PD-L1 transcript, H-score HALO})$ vs $\rho(\text{PD-L1 transcript, H-score QuPath})$	0.74	0.76	-0.43	0.67	No significant difference	
$\rho(\text{PD-L1 transcript, pix H-score VIS})$ vs $\rho(\text{PD-L1 transcript, pix H-score HALO})$	0.80	0.79	0.19	0.85	No significant difference	
$\rho(\text{PD-L1 transcript, pix H-score VIS})$ vs $\rho(\text{PD-L1 transcript, H-score HALO})$	0.80	0.74	1.16	0.26	No significant difference	
$\rho(\text{PD-L1 transcript, pix H-score VIS})$ vs $\rho(\text{PD-L1 transcript, H-score QuPath})$	0.80	0.76	0.80	0.43	No significant difference	
$\rho(\text{PD-L1 transcript, pix H-score VIS})$ vs $\rho(\text{PD-L1 transcript, ATM score})$	0.80	0.56	3.29	0.003	Significant difference	
$\rho(\text{PD-L1 transcript, pix H-score HALO})$ vs $\rho(\text{PD-L1 transcript, H-score HALO})$	0.79	0.74	1.00	0.33	No significant difference	
$\rho(\text{PD-L1 transcript, pix H-score HALO})$ vs $\rho(\text{PD-L1 transcript, H-score QuPath})$	0.79	0.76	0.47	0.64	No significant difference	
$\rho(\text{PD-L1 transcript, pix H-score HALO})$ vs $\rho(\text{PD-L1 transcript, ATM score})$	0.79	0.56	2.48	0.02	significant difference	

Table 2. Table lists the results of William's t-test to test for significant difference in the Spearman correlation coefficients between PD-L1 mRNA transcript or pathologist H-score and different DIA endpoints.

Statistical analysis of correlation coefficients for 5T4						
	ρ_{12}	ρ_{13}	z-score	p-value	Result	
Comparing pairwise correlations between DIA endpoint and Pathologist H-score						
$\rho(\text{Path H-score, H-score HALO})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.75	0.79	-0.49	0.63	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, pix H-score HALO})$	0.83	0.82	0.31	0.76	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, H-score HALO})$	0.83	0.75	1.69	0.11	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.83	0.79	0.49	0.63	No significant difference	
$\rho(\text{Path H-score, pix H-score VIS})$ vs $\rho(\text{Path H-score, ATM score})$	0.83	0.76	1.37	0.18	No significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, H-score HALO})$	0.82	0.75	1.46	0.16	No significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, H-score QuPath})$	0.82	0.79	0.37	0.71	No significant difference	
$\rho(\text{Path H-score, pix H-score HALO})$ vs $\rho(\text{Path H-score, ATM score})$	0.82	0.76	1.32	0.20	No significant difference	
Comparing pairwise correlations between DIA endpoint and 5T4 transcript						
$\rho(5T4 \text{ transcript, H-score HALO})$ vs $\rho(5T4 \text{ transcript, H-score QuPath})$	0.74	0.55	2.11	0.05	Significant difference	
$\rho(5T4 \text{ transcript, pix H-score VIS})$ vs $\rho(5T4 \text{ transcript, pix H-score HALO})$	0.76	0.74	0.67	0.51	No significant difference	
$\rho(5T4 \text{ transcript, pix H-score VIS})$ vs $\rho(5T4 \text{ transcript, H-score HALO})$	0.76	0.74	0.50	0.63	No significant difference	
$\rho(5T4 \text{ transcript, pix H-score VIS})$ vs $\rho(5T4 \text{ transcript, H-score QuPath})$	0.76	0.55	2.40	0.03	Significant difference	
$\rho(5T4 \text{ transcript, pix H-score VIS})$ vs $\rho(5T4 \text{ transcript, ATM score})$	0.76	0.69	1.40	0.18	No significant difference	
$\rho(5T4 \text{ transcript, pix H-score HALO})$ vs $\rho(5T4 \text{ transcript, H-score HALO})$	0.74	0.74	0.10	0.92	No significant difference	
$\rho(5T4 \text{ transcript, pix H-score HALO})$ vs $\rho(5T4 \text{ transcript, H-score QuPath})$	0.74	0.55	2.10	0.05	Significant difference	
$\rho(5T4 \text{ transcript, pix H-score HALO})$ vs $\rho(5T4 \text{ transcript, ATM score})$	0.74	0.69	1.11	0.28	No significant difference	

Table 3. Table lists the results of William's t-test to test for significant difference in the Spearman correlation coefficients between 5T4 mRNA transcript or pathologist H-score and different DIA endpoints.