

Analyst

Accepted Manuscript

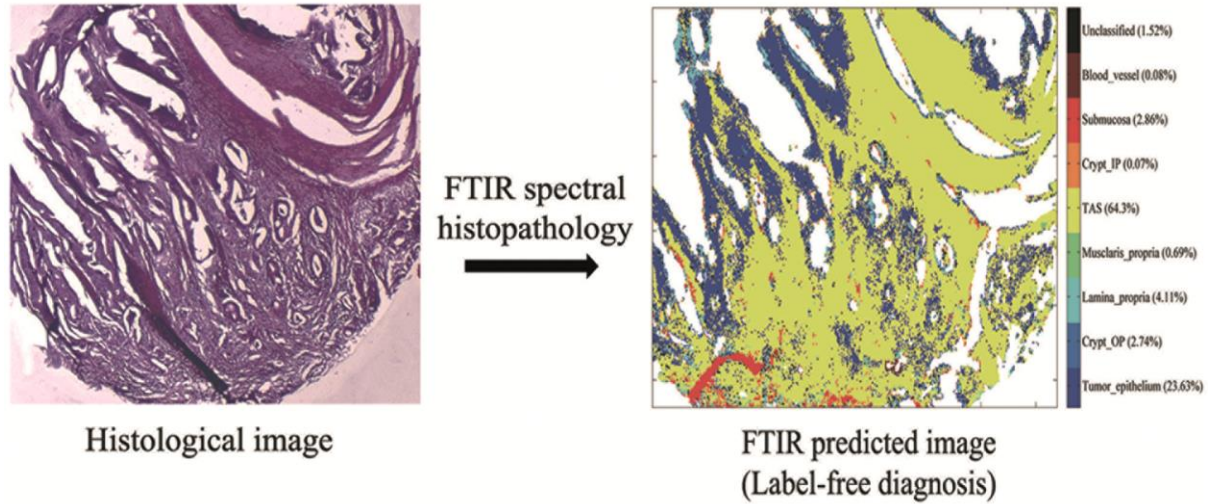


This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.



Automated and label-free colon cancer diagnosis and identification of tumor-associated features using FTIR spectral histopathology directly on paraffinized tissue arrays

1
2
3
4 1 **Title:**

5
6 2 Infrared spectral histopathology for cancer diagnosis; a novel approach for automated pattern recognition
7
8 3 of colon adenocarcinoma
9

10
11 4
12 5 **Authors and affiliations:**

13
14
15 6 Jayakrupakar Nallala^{a, b, e}, Marie-Danièle Diebold^{a, b, c}, Cyril Gobinet^{a, b}, Olivier Bouché^{a, b, d}, Ganesh
16
17 7 Dhruvananda Sockalingum^{a, b}, Olivier Piot^{a, b}, Michel Manfait^{a, b*}.
18
19

20
21 8 ^aUniversité de Reims Champagne-Ardenne, MéDIAN-Biophotonique et Technologies pour la Santé, UFR
22
23 9 de Pharmacie, 51 rue Cognacq-Jay, 51096 REIMS cedex, France.
24

25 10
26
27 11 ^bCNRS UMR7369, Matrice Extracellulaire et Dynamique Cellulaire, MEDyC, Reims, France.
28

29 12
30
31 13 ^cCHU de Reims, Laboratoire Central d'Anatomo-Pathologie, 51092 Reims Cedex, France.
32

33 14
34
35 15 ^dCHU de Reims, Service d'Hépatogastroentérologie et de Cancérologie Digestive, 51092 Reims Cedex,
36
37 16 France.
38

39 17
40
41 18 ^eCurrent Address: University of Exeter, Biomedical Physics, School of Physics, Exeter EX4 4QL, United
42
43 19 Kingdom.
44

45 20
46
47 21
48 22 **Correspondence:**

49 23 Michel Manfait*

50
51
52 24 michel.manfait@univ-reims.fr, Tel: +33 32 69 13 57 4, Fax: +33 32 69 13 55 0.
53
54
55
56
57
58
59
60

2

ABSTRACT:

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Histopathology remains the gold standard method for colon cancer diagnosis. Novel complementary approaches for molecular level diagnosis of the disease are need of the hour. Infrared (IR) imaging could be a promising candidate method as it probes the intrinsic chemical bonds present in a tissue, and provides a “spectral fingerprint” of the biochemical composition. To this end, IR spectral histopathology, which combines IR imaging and data processing techniques, was employed on seventy seven paraffinized colon tissue samples (48 tumoral, 29 non-tumoral) in the form of tissue arrays. To avoid chemical deparaffinization, a digital neutralization of the spectral interferences of paraffin was implemented. Clustering analysis was used to partition the spectra and construct pseudo-colored images, for assigning spectral clusters to various tissue structures (normal epithelium, malignant epithelium, connective tissue etc). Based on the clustering results, linear discriminant analysis was then used to construct a stringent prediction model which was applied on samples without a priori histopathological information. The predicted spectral images not only revealed common features representative of the colonic tissue biochemical make-up, but also highlighted additional features like tumor budding, tumor-stroma association in a label-free manner. This novel approach of IR spectral imaging on paraffinized tissues showed 100 % sensitivity and allowed detection and differentiation of normal and malignant colonic features based purely on their intrinsic biochemical features. This non-destructive methodology combined with multivariate statistical image analysis appears as a promising tool for colon cancer diagnosis and opens the way to the concept of numerical spectral histopathology.

3

1. INTRODUCTION:

Colorectal cancer has one of the highest incidence and mortality among all the cancers affecting both sexes, of which the type adenocarcinoma is the most common.¹ Radiation therapy, chemotherapy and surgical intervention have improved the life expectancy of cancer patients, but the outcome of these methods is dependent upon the stage and the accuracy in diagnosis.² Currently different detection and screening methods are employed for colorectal cancers, including fecal occult blood test (FOBT),³ sigmoidoscopy,⁴ colonoscopy,⁵ etc. However, the final diagnosis is settled upon the microscopic examination of the symptomatic tissue with the 'gold standard' histopathology in which preferential stains are used to enhance visualization of the tissue morphological alterations. Such alterations (pre-cancerous or cancerous) are the manifestations of the biomolecular changes that have already undergone the provocative changes for malignancy. However, the ongoing state of the tissue molecular changes during the onset or progression of malignancy, without any morphological signatures, poses a challenge for identification. In certain cases, immunohistochemistry (IHC) is used to identify specific proteins of interest which can give a molecular level understanding of the malignant condition. Histopathology requires precise human expertise which limits high-throughput diagnosis. Although, the histopathological diagnosis is based on morphological examination, it has successfully served in cancer diagnosis over several years. Additionally, if it is combined with approaches that could provide complementary biochemical information in a rapid, cost effective manner and reducing human involvement, the efficacy of the histopathological diagnosis can be completed.

In this regard, the optical spectroscopic approach of IR imaging appears as a potential candidate for routine tissue characterization, and has been exploited as a diagnostic tool on various tissues⁶⁻¹⁸ which also paved the way to the concept of spectral histopathology.¹⁹⁻²³ IR spectroscopy probes intrinsic chemical bond vibrations of biomolecules and thus provides a biochemical fingerprint of the tissues. Combined with an imaging set-up, spectral images can be obtained rapidly in a label-free manner, in which each pixel element harbors an IR spectrum containing biochemical information at each wavenumber. Such IR images can be exploited using computer based multivariate cluster analysis to generate digitally stained morphological maps of the tissue histology. Since the constituent IR spectra of each digitally stained histological class represent its biochemical signature, such as collagen features in

4

1
2
3 83 the connective tissue, specific spectral signatures can be identified from different histological classes.
4
5 84 Such signatures can be used to train predictive algorithms for identification of unknown tissues in a rapid
6
7 85 and user-friendly manner. One of the important possibilities of using this methodology is automation of
8
9 86 this protocol which can reduce human involvement and provide an objective biochemical based
10
11 87 diagnostic approach.

12
13 88 In this regard, we carried out spectral histopathology based on IR imaging in conjunction with multivariate
14
15 89 analysis. The main objectives were to digitally detect and identify malignancy and its associated features
16
17 90 on unknown tissues without any chemical staining, constituting an automated diagnosis for colon
18
19 91 adenocarcinoma. For this, 77 human colon tissues from normal and moderately differentiated
20
21 92 adenocarcinoma were analyzed, in the form of paraffinized tissue arrays that were stabilized in an
22
23 93 agarose matrix. The agarose matrix provides stability to the paraffinized tissue cores thereby reducing
24
25 94 tissue loss during microtome sectioning, and also facilitates handling of tissue array sections. The tissue
26
27 95 arrays are increasingly used in pathological studies since they constitute a large source of information
28
29 96 and permit high-throughput analysis for modern histological practices.²⁴ An innovative process of digital
30
31 97 deparaffinization was specially implemented to avoid chemical dewaxing, and also to reduce toxic
32
33 98 chemical treatments and time consumption.²⁰ Then, a prediction model representing the main colon
34
35 99 histological classes was constructed and its robustness was evaluated on subsequent number of tissue
36
37 100 array cores. Digital annotation using this model facilitated characterization of malignancy, and malignancy
38
39 101 associated features such as tumor budding, and tumor-stroma association.

102

103 **2. MATERIALS AND METHODS:**

104 **2.1. Sample preparation:**

105 Seventy seven formalin fixed paraffin embedded (FFPE) colon tissue samples (48 tumoral and 29 non-
106 tumoral) from 32 cancer patients were obtained from the Reims University Hospital, with the approval of
107 the Institutional Review Board. All the tumoral samples were moderately differentiated colon
108 adenocarcinoma with the TNM grade ranging from T3N0M0 to T4N2M0. The sample details are
109 presented in Supplementary Table 1. Several paraffinized tissue arrays that were stabilized in an agarose
110 matrix were manually prepared from these samples. A single sample spot in the tissue array block was

5

1
2
3 111 approximately 3 mm in diameter. For each tissue array consisting around 12-16 spots, 3 and 10 μm thick
4
5 112 sections (adjacent in most cases) were obtained. While the 3 μm section was used by the pathologist for
6
7 113 conventional histopathological analysis via hematoxylin, phloxine, and saffron (HPS) staining, the first 10
8
9 114 μm unstained section was used for IR imaging analysis and the second stained section for additional
10
11 115 histopathological comparison. The HPS stained sections were chemically deparaffinized while the
12
13 116 unstained tissue section for IR imaging was mounted on an IR compatible calcium fluoride (CaF_2) support
14
15 117 without any chemical deparaffinization.
16
17 118

19 119 **2.2. Instrumentation and FTIR data collection:**

20
21 120 IR images were acquired, by an IR imaging system (Spotlight 300, Perkin Elmer, Courtaboeuf, France)
22
23 121 equipped with liquid nitrogen-cooled 16-element MCT detector, at $6.25 \times 6.25 \mu\text{m}^2$ pixel size, and 4 cm^{-1}
24
25 122 spectral resolution averaged to 16 scans, in the mid-IR range of 750 to 4000 cm^{-1} . The system was
26
27 123 continuously purged with dry air. The background spectrum from the CaF_2 support was recorded each
28
29 124 time prior to image acquisition, using the same parameters as that of the IR image. The methodology for
30
31 125 FTIR spectral imaging of tissue arrays is represented in Supplementary Figure 1. A total of 8 141 566 IR
32
33 126 spectra were recorded from 77 images at an average of 105 734 per image owing to the large size of the
34
35 127 tissue array spots, and the high spatial resolution selected for imaging.
36
37 128

38 129 **2.3. Data pre-processing:**

39
40 130 Raw IR data was corrected from various spectral interferences. An atmospheric correction was performed
41
42 131 to remove contribution from water vapour and CO_2 by the built-in Perkin Elmer Spotlight software and
43
44 132 further processing was carried out using programmes written in Matlab 7.2 (The Mathworks, Natick, MA).
45
46 133 The spectra were reduced to the IR absorption range of $900\text{-}1800 \text{ cm}^{-1}$ that contains several informative
47
48 134 biochemical vibrations^{25,26} as far as the tissue features are considered. Neutralization of paraffin and
49
50 135 agarose contributions was carried out using a modified Extended Multiplicative Signal Correction (EMSC).
51
52 136 In addition to paraffin model, a correction model for agarose was inserted into the EMSC algorithm.^{20, 27}
53
54 137 As detailed, the EMSC algorithm neutralizes the influence of their spectral variabilities by a modeling
55
56 138 procedure rather than directly subtracting the spectral signatures of paraffin and agarose.²⁰ Therefore it is
57
58
59
60

6

1
2
3 139 important to note that paraffin and agarose features are not removed, but their spectral variabilities are
4
5 140 neutralized. Therefore, in the image analysis only the spectral variabilities originating from the
6
7 141 biochemical features are taken into account rather than those from physical features of paraffin and
8
9 142 agarose which are no longer apparent. Furthermore, EMSC has been adapted to address the inter-
10
11 143 patient variability using a single target spectrum (also called the model spectrum or the reference
12
13 144 spectrum) for the all the tissue samples. Using the same target spectrum for all the samples has been an
14
15 145 important criterion in our application in order to correct all the spectra from the same amount of baseline,
16
17 146 paraffin and agarose, while keeping the biochemical information specific to each sample. If a different
18
19 147 target spectrum is used for each sample (e.g. for normal and cancerous tissue), the corrected spectra of
20
21 148 each sample will have different shapes mainly because of the different baselines and/or paraffin signals
22
23 149 and/or agarose signals composing each of the target spectra, and not because of the biochemical
24
25 150 differences between normal and cancer tissues.²⁸⁻²⁹ The IR spectra were also corrected for baseline and
26
27 151 then normalized using the same algorithm. Outliers (N=3 335 684 spectra) in the form of paraffin and
28
29 152 agarose spectra, and spectra with poor signal-noise ratio were eliminated from the analysis and were
30
31 153 depicted as white pixels in all the IR images.

32
33 154

34 35 155 **2.4. Data processing:**

36
37 156 The pre-processed data (N=4 805 882 spectra) was subjected to multivariate statistical prediction
38
39 157 analysis. For this, spectral data from the non-tumoral and the tumoral samples was separated into a
40
41 158 training group (N=9, Supplementary Table 1, sample # TG), and a validation group (N=68). While the
42
43 159 training group, representing the IR spectral signatures indicative of malignancy and other histological
44
45 160 structures, was used for construction of a prediction model based on linear discriminant analysis (LDA),
46
47 161 the validation group (external validation) was used for validating the model on unknown samples for
48
49 162 automatic recognition of tissue features, to enable identification of malignancy. LDA is a multivariate
50
51 163 supervised statistical technique that aims at maximizing the between-class variance and minimizing the
52
53 164 within-class variance and has been exploited in various studies.^{26-27, 30-31}

54
55 16556
57 166

58

59

60

7

167 **2.4.1. Cluster analysis for LDA training:**

168 The huge number of IR spectra from each image corresponding to the training group was subjected to
169 unsupervised k-means clustering method owing to its capability of rapid and huge data clustering.³² This
170 method iteratively partitions the spectra into different clusters based on the spectral signatures from the
171 intrinsic biochemical composition of the tissue. Therefore, spectra with similar biochemical characteristics
172 group into the same cluster. In k-means clustering, each spectrum belongs to a unique cluster and can
173 thus be represented by one color. K-means clustering performed using defined cluster numbers resulted
174 in the construction of digital color-coded images. These were then compared to adjacent HPS stained
175 sections to annotate by an expert pathologist, each spectral cluster to the tissue structural feature that it
176 corresponds to. The spectral distance between different k-means clusters was visualized in a dendrogram
177 obtained by hierarchical clustering analysis using Ward's linkage algorithm.

179 **2.4.2. Prediction model:**

180 The initially k-means clustered and annotated spectra were used as inputs for the LDA model. Training
181 group spectra (Supplementary Table 1 # TG) from 9 samples across 6 different patients were considered
182 for the model, to take into account the inter-patient variability. The prediction model consisted of 8 classes
183 with different number of spectra, representing various histological features of non-tumoral and tumoral
184 tissues: the normal epithelium defined by the crypt inner-part (Crypt-IP) (N = 8377) and the crypt outer-
185 part (Crypt-OP) (N = 3567), the lamina propria (N = 14 106), the submucosa (N = 3964), the tumor
186 epithelium (N = 35 083), the tumor-associated stroma (N = 16 409), the blood vessel (N = 782) and the
187 muscularis propria (N = 4514). These spectra (N=86 802) constituting one-third of the spectra from each
188 class were used to train the model and the other two-thirds were used for an internal validation to
189 optimize the model. The prediction model was then applied in an external validation on different unknown
190 samples, the spectra from which were secluded from the model, to evaluate its robustness. The external
191 validation consisted of 68 samples encompassing a large scale spectral data base of 4 130 879 spectra.
192 It has to be noted that if only the number of patients used in the external validation was to be considered
193 (instead of the number of samples from all the patients as is the case in this study) the external validation
194 group consisted of 26 patients, since several samples were obtained from a single patient

8

1
2
3 195 (Supplementary Table 1). The predictions were carried out in the IR spectral range of 1080 cm^{-1} - 1300
4
5 196 cm^{-1} , at a posterior probability of 0.5, wherein for each pixel a probability of belonging to each class is
6
7 197 calculated, and the pixel showing the highest probability is assigned to a class. If the highest probably is
8
9 198 inferior to the posterior probability of 0.5, the pixel is termed as unclassified and is not attributed to any
10
11 199 class. The final model based diagnosis of cancer by the presence of tumor pixels was confirmed by the
12
13 200 presence of tumoral areas in the corresponding region of the HPS stained tissue, using the gold-standard
14
15 201 histopathological validation. Validation based on the presence of certain number of pixels (tumor pixels)
16
17 202 was not considered as a dedicated approach in this study where heterogeneous tumoral tissue types are
18
19 203 considered which contain varying amount of tumoral cells.
20
21 204

22 205 **2.4.3. Spectral information to biochemical information (spectral analysis):**

23 206 Since the spectral signatures are based on the biochemical properties of the tissue features, it was
24
25 207 attempted to characterize the biochemical alterations characteristic of malignancy and the relationship of
26
27 208 malignant tissue with the surrounding stroma. For this, the Mann-Whitney *U* test was applied to compare
28
29 209 spectra from selected cluster groups used in the prediction model training in order to identify the most
30
31 210 discriminant wavenumbers.
32
33 211

34 211

35 212 **2.5. Immunohistochemistry (IHC):**

36 213 IHC was used as a complementary tool (on adjacent sections) to enhance visibility of tumor budding
37
38 214 (Anti-Human Cytokeratins-large spectrum Monoclonal Antibody, Clone KL 1, dilution 1/50, Immunotech,
39
40 215 France) and to precise the nature of the inflammatory cells: T-lymphocytes (CD3 Rabbit anti-Human
41
42 216 Polyclonal Antibody, dilution 1/200, Dako, France), and B-lymphocytes (CD20 Mouse antibody, clone L6
43
44 217 mouse, dilution 1/400, Dako, France), in order to validate some of the important observations detected by
45
46 218 IR spectral imaging. This was performed using the fully automated IHC staining protocol (XT ultraView
47
48 219 DAB v3).
49
50 220

51 221

52 222

53

54

55

56

57

9

3. RESULTS:**3.1. Cluster analysis:**

K-means clustering was used to identify the spectral signatures characteristic of the main histological features of the non-tumoral and the tumoral colon tissues, which permitted construction of digitally stained images. For the non-tumoral as well as the tumoral tissues, this approach permitted to identify, and to recover automatically the important histological components in comparison to the adjacent HPS stained images as shown in the figure 1 (Supplementary Table 1, sample # 1D and 12C). As an example, for the non-tumoral colon tissue (figure 1A) 8 clusters permitted the observation of the important histological structures representing the colon tissue organization. They included the colon mucosa constituted by well-differentiated crypts (cluster 8 - inner part and cluster 6 - outer part); and the lamina propria (cluster 1), the supportive loose connective tissue in which the crypts are organized. The residual mucin (cluster 2) was observed to be localized within the crypt lumen while a small amount was seen secreted outside. The submucosa, attributed to clusters 4, 5 and 7 was distinguished effectively from the lamina propria by the clustering method. Finally cluster 3 appeared to represent the blood vessels. On the contrary, in the typical adenocarcinomatous tissue (figure 1B), the only important histological classes retrieved were the tumor epithelium (cluster 1) and its associated stroma in the tumor vicinity (cluster 6). Most of the other clusters represented the fibrous stromal tissue. The corresponding dendrogram showed the close spectral nature of the tumor associated stroma to its tumor where they are very closely grouped (clusters 1 and 6) while the stroma that is not in direct contact with the tumor epithelium appear more distant. A total of 11 clusters were required to identify these features. In both cases, considering the overall colon tissue organization, increasing the number of clusters did not add any further retrievable histological information. The k-means clustering is an efficient method to identify IR spectral markers specific to different histological components of non-tumoral and tumoral colon tissues. On the basis of these spectral signatures, the diagnostic potential of IR spectral imaging has been evaluated using a LDA based prediction model as schematically represented in Supplementary Figure 2.

3.2. Optimization of the prediction model - internal validation group:

10

1
2
3 250 The LDA based prediction model developed from 9 samples (6 patients) with 8 different classes
4
5 251 comprising a total of 86802 spectra was trained, and tested in an internal validation. The sensitivity of the
6
7 252 prediction model in the internal validation can be evaluated from the confusion matrix which shows the
8
9 253 agreement between the histopathological class annotation (real class) and the IR spectral prediction
10
11 254 (predicted class) (Table 1). Different spectral regions were tested and the highest sensitivity (average
12
13 255 89.38%) was obtained for the region between 1080 cm^{-1} to 1300 cm^{-1} . It has to be noted that for the class
14
15 256 tumor epithelium a specificity of 96.4 % was reached, and showed no confusion with the class normal
16
17 257 epithelium (comprising crypt inner and outer parts).
18
19 258

20 21 259 **3.3. Tumor detection and tissue characterization in unknown samples - external validation group:**

22
23 260 The external validation was performed on the remaining 68 blind samples involving a large scale spectral
24
25 261 bank of 4 130 879 spectra and showed 100 % sensitivity for the tumor class. Along with tumor class,
26
27 262 other histological classes were also identified with high correlation to the conventional histology.

28
29 263 A representative demonstration of prediction on unknown non-tumoral and tumoral samples is shown in
30
31 264 figure 2 (Supplementary Table 1, sample # 14D and 7C). The figure 2A histologically corresponded to a
32
33 265 non-tumoral colon tissue in which the prediction model correctly identified its characteristic features with
34
35 266 similar morphological attributes to that of the histological image. Counterpart to the normal tissue,
36
37 267 histologically the figure 2B corresponded to a typical moderately differentiated colon adenocarcinoma. In
38
39 268 this, the spectral characteristics of the normal mucosa were absent and the only distinguished ones were
40
41 269 malignant epithelial component with its associated stroma. Additionally, identification of features difficult
42
43 270 to discern using conventional techniques, such as tumor budding was facilitated.
44
45 271

46 272 **3.4. Detection and characterization of malignancy associated features:**

47 48 273 **3.4.1. Tumor budding:**

49
50 274 Budding is characterized by small clusters of isolated tumor cells which become detached from the
51
52 275 neoplastic epithelium and migrate into the stroma, and is an indication of high tumor invasiveness in
53
54 276 colorectal cancers. Although this morphological phenomenon is detectable in conventional histopathology
55
56 277 at high power magnification, IHC may be employed for better visualization. The IR prediction model was
57
58
59
60

11

1
2
3 278 able to clearly identify this tumor particularity even in the presence of abundant stroma as shown in the
4
5 279 figure 3, (Supplementary Table 1, sample # 9B). In the same tumoral sample, along with the malignant
6
7 280 epithelium, there was presence of some normal epithelial component together with normal connective
8
9 281 tissue, and all these features were identified by the prediction model. Importantly, both the malignant and
10
11 282 the non-malignant epithelial cells were selectively stained and discriminated using a specific color-code.
12
13 283 The positive staining of the epithelial cells can be seen in the IHC image (see figure 3, right panel).
14
15 284 Another tissue section obtained from different position (Supplementary Table 1, sample # 9A) of the same
16
17 285 tumor also showed tumor budding in a stroma dominant environment, and each time it was identified by
18
19 286 the prediction model, which was later confirmed by IHC studies (Supplementary Figure 3).
20
21 287

22 23 288 **3.4.2. Tumor stroma association:**

24
25 289 The tumor-stroma association was also reported using IR spectral imaging. The confusion matrix (table 1)
26
27 290 highlighted the spectral proximity of tumor and its associated stroma in which, indeed 16.3 % of tumor
28
29 291 associated stroma pixels were classified in the tumor class. Complementarily, in the predicted images
30
31 292 these two classes appeared in geographic proximity (figure 4) (Supplementary Table 1, sample # 11B). In
32
33 293 the same image, distinction between the tumor associated stroma and the normal connective tissue
34
35 294 corresponding to the submucosa was attained, while in the histological stained section, this was
36
37 295 indistinguishable. The above mentioned tumor-stroma features were also observed in the other tumoral
38
39 296 samples (Supplementary Table 1, sample # 11A, 11C, 12A, 13A, and 15A) as shown in Supplementary
40
41 297 Figure 4 including the cases of budding (fig 3).
42
43 298

44 45 299 **3.5. Vibrational analysis of spectroscopic markers:**

46
47 300 In this study, the k-means clustering was performed using the IR spectral range of 900 cm^{-1} - 1800 cm^{-1}
48
49 301 that enabled identification and attribution of the important colon histological classes. For unknown sample
50
51 302 prediction, this zone was narrowed down to 1080 cm^{-1} to 1300 cm^{-1} harboring some of the important
52
53 303 biomolecular vibrational modes implicated in colon cancers, and which showed the best prediction
54
55 304 outcome for all the classes together. As shown in figure 5, the most discriminant wavenumbers within this
56
57 305 zone were identified by the Mann-Whitney U test performed on the individual spectra and represented on
58
59
60

12

1
2
3 306 the average spectra for the following pair-wise comparisons: normal epithelium with malignant epithelium
4
5 307 (adenocarcinoma) for understanding the molecular alterations characteristic of malignancy;
6
7 308 adenocarcinoma with its associated stroma to understand the tumor induced alterations in the stromal
8
9 309 tissue; and the normal connective tissue with the tumor associated stroma. From the discriminant
10
11 310 wavenumbers identified for all comparisons, a tentative correlation of IR vibrations to the biomolecular
12
13 311 information was attempted as shown in Supplementary Table 2. Importantly, comparing the normal
14
15 312 epithelium with the tumoral epithelium, the main differences in the IR peaks were attributed to symmetric
16
17 313 and asymmetric PO_2^- vibrations of the nucleic acids that demonstrated relatively higher intensities in
18
19 314 normal than the tumoral tissues. Similarly, the C-O stretching vibration corresponding to carbohydrates
20
21 315 was relatively more intense in normal than the tumoral tissues. At the same time the hydrogen bonded C-
22
23 316 O groups of proteins in the normal epithelium was observed to be decreased in the tumoral epithelium,
24
25 317 while the opposite tendency was observed for the non-hydrogen bonded C-O groups of proteins.
26
27 318 Secondly, when comparing adenocarcinoma with tumor associated stroma, and tumor associated stroma
28
29 319 with connective tissue, the discriminating spectral features appeared to be contributed principally from
30
31 320 collagen features.

32
33 321

34 322 **4. DISCUSSION:**

35
36 323 Spectral histopathology based on IR imaging has been carried out to develop an innovative label-free
37
38 324 diagnostic methodology directly on FFPE tissue arrays embedded in an agarose matrix without any
39
40 325 chemical pre-treatments. EMSC that has been initially developed to separate light scattering effects from
41
42 326 light absorbance effects, has also been used for accomplishing neutralization of paraffin contributions in
43
44 327 IR spectral analysis.^{9,27,33-35} In this study, both paraffin and agarose interferences on the IR spectral
45
46 328 images have been neutralized digitally without the use of any chemicals, using an improved EMSC
47
48 329 algorithm. One of the important advantages of using of paraffinized tissues stabilized in an agarose matrix
49
50 330 is that the scattering effects such as Mie scattering due to the differences in the refractive indices of the
51
52 331 media are reduced by index matching. Additionally, resonant Mie-scattering that is related to a physical
53
54 332 phenomenon and which can cause peak shape distortion and peak shift (e.g., the amide I peak), resulting
55
56 333 in unreliable chemical interpretation is also reduced.³⁶

13

4.1. Clustering:

K-means clustering provided a rapid way to classify the IR spectral images into their constituent histological classes in comparison to the chemically stained conventional images. While the non-tumoral colon tissues were characterized by well-differentiated architecture with both inner and the outer cryptal parts clearly distinguishable together with the connective tissue, the malignant tissues which were all of the advanced colon cancer types, were characterized by the loss of differentiation of the normal colon glands with no visible lumen; and presence of stromal tissue. The digital staining of each k-means cluster formed the basis for spectral marker assignment comprising the malignant colon characteristics, along with the normal tissue features, at different organizational levels of the colon wall. Based on this spectral database from as little as 12 % of the samples, a prediction model was trained for automatic detection of malignancy in unknown specimens independently of conventional histopathology.

4.2. Prediction:

Some of the earlier IR imaging studies have tested prediction algorithms on different tissue types.^{25,26} However, the number of spectra used for constructing the model was limited compromising the robustness of the model. In our study, the relatively high resolution image acquisition parameters applied to tissue arrays (3 mm diameter) constituted a huge bank of 86 802 spectra in the prediction model, representative of the biochemical signatures of distinct colon structures, making it highly robust. Only one such IR imaging study on prostate tissues has used such a robust model for prediction on unknown tissues.³⁷ In this study, 8 classes were included that described the colon tissue organization in non-tumoral and tumoral samples. Even with a high sensitivity of the model (such as in the case of tumor budding), some of these histological structures may share certain similar molecular constituents with other histological classes present in the model (tumor and tumor associated stroma), or not present in the model (muscularis mucosa and tumor associated stroma). The spectral proximity arising from this leads to misclassification between such classes as shown in the Supplementary Figure 5, concerning the muscularis mucosa (visible in the HPS image) which is identified as tumor associated stroma (Supplementary Table 1, sample # 27). It has to be noted that there was no class for the muscularis mucosa in the model. This attribution can be presumed to have arisen from the residual normal

14

1
2
3 362 muscularis mucosa signatures present in the tumor associated stroma from which the corresponding
4
5 363 class was constructed in the prediction model. This prediction error appeared predominantly in non-
6
7 364 tumoral samples where there is an intact muscularis mucosa. Despite these misclassifications, an overall
8
9 365 high correlation between the predicted spectral classes and the corresponding histological structures is
10
11 366 observed in the confusion matrix.

12
13 367

15 368 **4.3. External validation:**

16
17 369 The remaining 88 % of IR spectral images were identified by the prediction model without any a priori
18
19 370 knowledge on their histopathology (external validation). These blind samples constituted a huge number
20
21 371 of 4 130 879 spectra that were scanned and annotated by the automated computer trained prediction
22
23 372 algorithm. The diagnosis was confirmed by an expert pathologist by using the conventional histological
24
25 373 images based on which a 100 % accuracy of the prediction model was obtained for tumor diagnosis. This
26
27 374 high sensitivity after scanning such a huge number of unknown spectra signifies the potential of the
28
29 375 current methodology as a diagnostic tool. The prediction analysis also facilitated simultaneously some
30
31 376 important malignancy associated features.

32
33 377

35 378 **4.4. Tumor budding:**

36
37 379 The phenomenon of tumor budding is of crucial clinical importance in colorectal cancers since it has been
38
39 380 shown to be a strong adverse prognostic marker.³⁸ As such, studies have correlated its occurrence with
40
41 381 aggressiveness and lymph node metastasis.³⁹ In this study, the prediction model facilitated the
42
43 382 identification of tumor budding in a stroma-dominant environment in an automated manner. This rapid
44
45 383 and selective detection of small clusters of isolated tumor cells in an abundant stroma environment
46
47 384 demonstrates the sensitivity and the applicability of the methodology avoiding the need of any histological
48
49 385 or immunological markers. This envisages an important prospect since the tumor de-differentiation in the
50
51 386 form of budding is being acknowledged as a key component in the metastatic process even in well- and
52
53 387 moderately differentiated tumors.^{40,41} At the same time, the color code based selective staining of the
54
55 388 epithelial counter parts in the same tissue shows the discriminatory ability and the biomolecular specificity
56
57 389 of this methodology.

58
59
60

15

4.5. Spectral Analysis:

The IR spectral region from 1000 cm^{-1} to 1300 cm^{-1} has been reported to carry important biochemical vibrations implicated in colon cancers and have been used for differentiating the malignant tissues from their normal counterparts.^{42,43} In this study, the most discriminant spectral wavenumbers were associated with relatively decreased intensities of symmetric and asymmetric PO_2^- vibrations of the nucleic acids in the tumoral epithelium when compared to the non-tumoral tissues. On contrary to the expected increased nucleic acid intensities as shown in several studies, these spectral changes corresponding to the biochemical alterations corroborate with some of the previous studies on colon cancers where the nucleic acid intensities were shown to be reduced in malignant conditions.^{32,44} It may be likely that the spectral changes involving nucleic acids are small in moderately differentiated tumors when compared to normal colon epithelial cells which themselves are highly proliferative in nature. One study has stated that decreased phosphate content in malignant colon tissues may be due to decrease in carbohydrate content,⁴⁵ which in our study was also indicated by the relatively less intense C-O stretching vibration corresponding to carbohydrates in the tumoral tissue than the normal. At the same time, the relative intensities of H-bonded C-O vibrations of proteins were observed to be more pronounced in the normal epithelium than the tumoral, while the non-H-bonded C-O bond vibrations were more pronounced in the tumor. These changes may be indicative of the molecular alterations associated with the amino acid side chains concerning tyrosine, serine and threonine.^{2,32,45,46} The molecular changes involving adenocarcinoma and tumor-associated stroma, and tumor associated stroma with connective tissue appear principally due to collagen features.

410

4.6. Tissue inflammation influences the model specificity:

In 12 out of 29 samples histologically described as non-tumoral (Supplementary Table 1, sample # LF); tumoral characteristics (over 4 % of pixels) were observed either regionally clustered or dispersed in the lamina propria, showing a specificity of 59%. The HPS images gave insight into the regionally clustered tumor pixels as corresponding to lymphoid follicles in the colon tissue. These structures showed spectral signatures close to the tumor group relative to the other classes. However, the tumor pixels dispersed in the lamina propria could not be accounted for as no visible correspondence between them and any

16

1
2
3 418 histological feature could be found in the HPS images. Since these tissues showed high inflammatory
4
5 419 infiltration, immuno-staining for T-lymphocytes (CD 3), B-lymphocytes (CD 20) and macrophages (KP 1)
6
7 420 was performed to verify if the dispersed pixels corresponded to the inflammatory cells. The positive
8
9 421 staining indicated that these pixels indeed corresponded mainly to interstitial T-lymphocytes as
10
11 422 representatively shown in the figure 6A (Supplementary Table 1, sample # 32). In parallel, the B-
12
13 423 lymphocytes were seen assembled in lymph follicles. Non-tumoral tissues without any marked
14
15 424 inflammation as confirmed by the IHC showed no tumor pixels in the IR spectral images (figure 6B)
16
17 425 (Supplementary Table 1, sample # 31). Since the model did not take into account inflammatory conditions
18
19 426 (because of the tissue complexity arising from polymorphisms of the inflammatory infiltrates in colon
20
21 427 cancers: polymorph predominant, mononuclear predominant, mixed or rich in lymphoid follicles, and the
22
23 428 difficulty to have a representative spectral signature), these features were attributed to the spectrally
24
25 429 nearest class which turned out to be the tumor class.

26
27 430 A recent IR imaging study on cervical cancer tissues also quoted the influence of inflammatory signatures
28
29 431 on the prediction model sensitivity and specificity.²³ To have a broader insight into this aspect, we further
30
31 432 looked at the spectral class attribution threshold for the tumor class. It turned out that for the attribution of
32
33 433 spectra to tumor class, majority of the spectra corresponding to the inflammatory signatures have lesser
34
35 434 threshold values compared to the tumor in which the majority of the spectra have the highest posterior
36
37 435 probability values (Supplementary Figure 6). Altogether, the IR signatures from the inflammatory regions
38
39 436 appeared to class spectrally closer to tumor than other classes of the prediction model indicating an
40
41 437 intermediate stage between normal and malignant condition, as was shown in an earlier study.⁴⁷

42
43 438 The current work of IR spectral imaging on colon tissues provides automated diagnosis of malignancy on
44
45 439 unknown samples. Various diagnostic features associated with malignancy which provides
46
47 440 complementary information are also characterized. Important features such as tumor budding, tumor-
48
49 441 stroma association are dealt with in a non-destructive and label-free manner. The analysis of such a large
50
51 442 spectral database makes the study all the more representative. All these features have never been dealt
52
53 443 together in colon cancer diagnosis using IR spectral imaging of paraffinized tissues in any of the previous
54
55 444 studies. IR spectral imaging presents an optimistic overture for cancer knowledge in modern
56
57 445 histopathology.
58
59
60

17

1
2
3 446 The current prediction model representing the important histological features of a colon tissue certainly
4
5 447 holds aspects for amelioration. The spectral attribution identified the inflammatory signatures classed
6
7 448 close to the tumor. Since these specific biochemical signatures were picked up by the model, the
8
9 449 inflammatory infiltration, which pose risk of developing into cancers, could be incorporated into the model
10
11 450 for an automated evaluation and direct diagnostic approach for inflammatory diseases. In the same
12
13 451 manner, classes' specific to early neoplastic condition such as dysplasia could be incorporated into the
14
15 452 model and their spectral attribution thresholds compared to that of adenocarcinoma and normal
16
17 453 epithelium. This can potentially provide insights, into spectral alterations in early neoplastic conditions and
18
19 454 therefore, for early diagnosis of cancers. Aspects like genotype specific tumoral signatures and their
20
21 455 treatment response sensibility unknown till now could open a new additional classification. Further, an
22
23 456 automated quantification can be achieved for features like amount of tumor presence, or the amount of
24
25 457 tumor budding, only limit being the use of adjacent tissue sections which may present slight variations
26
27 458 from the reference tissue.

28
29 459

30 31 460 **5. CONCLUSIONS:**

32
33 461 The IR spectral imaging combined with multivariate statistical analyses appears as an optimistic
34
35 462 diagnostic approach for colon cancers in complement to conventional histopathology. This innovative
36
37 463 imaging approach enabled direct analysis of paraffinized tissue arrays and, via the employment of
38
39 464 mathematical deparaffinization the need for chemical pretreatments was reduced. The prediction model
40
41 465 permitted identification of unknown samples with a very high sensitivity, while the false positive prediction
42
43 466 in the non-tumoral samples has put forth the influence of the inflammatory component. This very large
44
45 467 scale spectral data base analyzed both in terms of training and validation shows the potentials of the IR
46
47 468 spectral imaging methodology for automated diagnostic purposes. Moreover, it eliminated the need for
48
49 469 sample staining and a priori knowledge of the sample to be analyzed. These optimistic results open a
50
51 470 new way for developing spectral biomarkers and libraries which could be used, in complement to
52
53 471 conventional histopathology, for early diagnosis, and also potentially for prognosis and theranostics of
54
55 472 cancers.

56 473
57
58
59
60

18

474 **ACKNOWLEDGEMENTS:**

475 This study was supported by a grant of Institut National du Cancer (INCa) and Canceropôle Grand Est.
476 We would like to thank Ligue contre le Cancer, Conférence de Coordination Interrégionale du Grand-Est,
477 and CNRS Projets Exploratoires Pluridisciplinaires, for financial support. Plateforme IBiSA "Imagerie
478 Cellulaire et Tissulaire", and the Tumorotheque, Champagne-Ardenne is also acknowledged. NJ is a
479 recipient of doctoral fellowship from the Région Champagne-Ardenne.

480

481 **REFERENCES:**

482

483 1. J. Ferlay, H. R. Shin, F. Bray, D. Forman, C. Mathers and D. M. Parkin, *Int. J. Cancer*, 2010, **127**,
484 2893-2917.

485

486 2. C. Conti, P. Ferraris, E. Giorgini, C. Rubini, S. Sabbatini, G. Tosi, J. Anastassopoulou, P. Arapantoni,
487 E. Boukaki, S. Konstadoudakis, T. Theophanides and C. Valavanis, *J. Mol. Struct.*, 2008, **881**, 46-51.

488

489 3. H. Miyoshi, M. Oka, K. Sugi, O. Saitoh, K. Katsu and K. Uchida, *Intern. Med.*, 2000, **39**, 701-706.

490

491 4. T. J. Zuber, *Am. Fam. Physician*, 2001, **63**, 1375-1380.

492

493 5. D. K. Rex, *Colon tumors and colonoscopy*, 2000, **32**, 874-883.

494

495 6. A. Tfayli, O. Piot, A. Durlach, P. Bernard and M. Manfait, *Biochim. Biophys. Acta*, 2005, **1724**, 262-269.

496

497 7. H. Fabian, N. A. Thi, M. Eiden, P. Lasch, J. Schmitt and D. Naumann, *Biochim. Biophys. Acta*, 2006,
498 **1758**, 874-882.

499

500

501

502

503

19

- 1
2
3 500 8. W. Steller, J. Einenkel, L. C. Horn, U. D. Braumann, H. Binder, R. Salzer and C. Krafft, *Anal. Bioanal.*
4
5 501 *Chem*, 2006, **384**, 45-154.
6
7 502
8
9 503 9. A. Travo, O. Piot, R. Wolthuis, C. Gobinet, M. Manfait, J. Bara, M. E. Forgue-Lafitte and P.
10
11 504 Jeannesson, *Histopathology*, 2010, **56**, 921-931.
12
13 505
14
15 506 10. M. J. Nasse, M. J. Walsh, E. C. Mattson, R. Reininger, A. Kajdacsy-Balla, V. Macias, R. Bhargava
16
17 507 and C. Hirschmugl, *Nat. Methods*, 2011, **8**, 413-416.
18
19 508
20
21 509 11. M. J. German, A. Hammiche, N. Ragavan, M. Tobin, L. J. Cooper, S. S. Matanhelia, A. C. Hindley, C.
22
23 510 M. Nicholson, N. J. Fullwood, H. M. Pollock and F. L. Martin, *Biophys. J*, 2006, **90**, 3783-3795.
24
25 511
26
27 512 12. K. Yano, S. Ohoshima, Y. Gotou, K. Kumaido, T. Moriguchi and H. Katayama, *Anal. Biochem*, 2000,
28
29 513 **287**, 218-225.
30
31 514
32
33 515 13. T. D. Wang, G. Triadafilopoulos, J. M. Crawford, L. R. Dixon, T. Bhandari, P. Sahbaie, S. Friedland,
34
35 516 R. Soetikno and C. H. Contag, *Proc. Natl. Acad. Sci. USA*, 2007, **104**, 15864-15869.
36
37 517
38
39 518 14. X. Zhang, Y. Xu, Y. Zhang, L. Wang, C. Hou, X. Zhou, X. Ling and Z. Xu, *J. Surg. Res*, 2011, **171**,
40
41 519 650-6.
42
43 520
44
45 521 15. C. Krafft, S. B. Sobottka, K. D. Geiger, G. Schackert and R. Salzer, *Anal. Bioanal. Chem*, 2007, **387**,
46
47 522 1669-1677.
48
49 523
50
51 524 16. J. Nallala, O. Piot, M. D. Diebold, C. Gobinet, O. Bouche', M. Manfait and G. D. Sockalingum,
52
53 525 *Cytometry Part A*, 2013, **83**, 294-300.
54
55 526
56
57 527 17. J. T. Kwak, S. M. Hewitt, S. Sinha, R. Bhargava, *BMC Cancer*, 2011, **11**, 62.
58
59
60

20

- 1
2
3 528 18. Rohit Bhargava. *Anal Bioanal Chem*, 2007, **389**, 1155-1169.
4
5 529
6
7 530 19. B. Bird, M. Miljkovic, S. Remiszewski, A. Akalin, M. Kon and M. Diem, *Lab. Invest*, 2012, 1-16.
8
9 531
10
11 532 20. J. Nallala, C. Gobinet, M. D. Diebold, V. Untereiner, O. Bouché, M. Manfait, G. D. Sockalingum and
12 533 O. Piot, *J. Biomed. Opt*, 2012, **17**, 1-12.
13
14 534
15
16 535 21. M. J. Walsh, S. E. Holtona, A. Kajdacsy-Ballab and R. Bhargava, *Vibrational Spectroscopy*, 2012, **60**,
17 536 23-28.
18
19 537
20
21 538 22. J. D. Pallua, C. Pezzei, B. Zelger, G. Schaefer, L.K. Bittner, V. A. Huck-Pezzei, S. A.
22 539 Schoenbichler, H. Hahn, A. Kloss-Brandstaetter, F. Kloss, G. K. Bonn and C. W. Huck, *Analyst*, 2012,
23 540 **137**, 3965-3974.
24
25 541
26
27 542 23. J. Eienkel, U. D. Braumann, W. Steller, H. Binder and L. Horn, *J. Histopathology*, 2012, **60**, 1084-
28 543 1098.
29 544
30
31 545 24. J. Kononen, L. Bubendorf, A. Kallionimeni, M. Bärlund, P. Schraml, S. Leighton, J. Torhorst, M.
32 546 Mihatsch, G. Sauter and O. Kallioniemi, *Nat. Med*, 1998, **4**, 844-847.
33 547
34
35 548 25. M. Khanmohammadi, A. B. Garmarudi, K. Ghasemi, H. K. Jaliseh and A. Kaviani, *Med. Oncol*, 2009,
36 549 **26**, 292-297.
37 550
38
39 551 26. M. Khanmohammadi, A. B. Garmarudi, S. Samani, K. Ghasemi and A. Ashuri, *Pathol. Oncol. Res*,
40 552 2010, **17**, 435-441.
41 553
42
43 554 27. E. Ly, O. Piot, R. Wolthuis, A. Durlach, P. Bernard and M. Manfait, *Analyst*, 2008, **133**, 197-205.
44 555
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

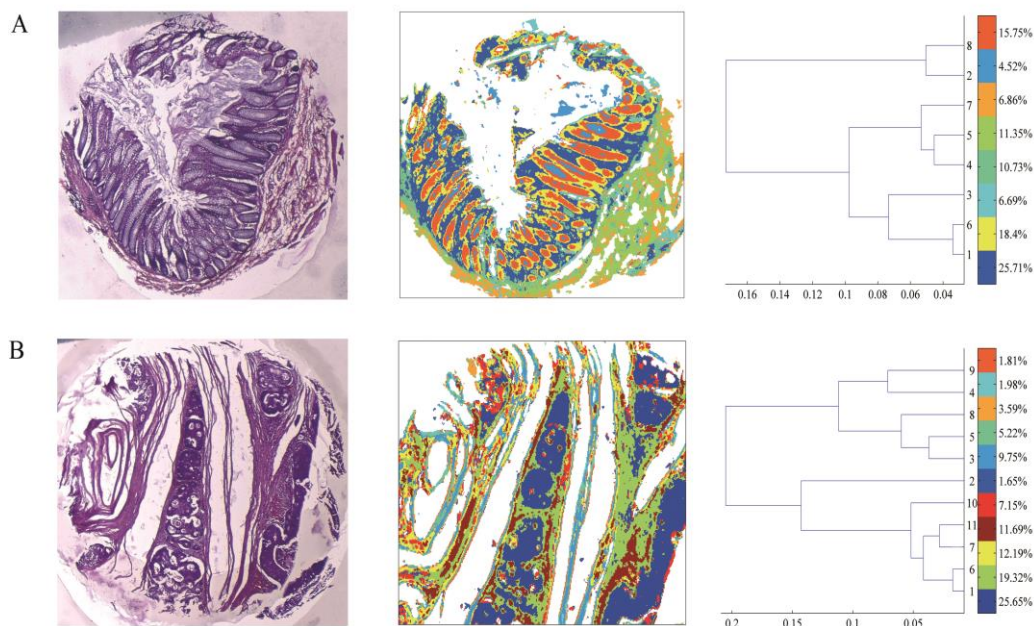
21

- 1
2
3 556 28. A. Kohler, N. K. Afseth and H. Martens, in Applications of Vibrational Spectroscopy in Food Science,
4
5 557 ed. E. Li-Chan, P. R. Griffiths and J. M. Chalmers, John Wiley & Sons, Ltd. 2010, vol. 1, pp. 89-97.
6
7 558
8
9 559 29. N. K Afseth and A. Kohler. *Chemometrics and Intelligent Laboratory Systems*, 2012, **117**, 92-97.
10
11 560
12
13 561 30. M. Khanmohammadi, M. A. Ansari, A. B. Garmarudi, G. Hassanzadeh and G. Garoosi, *Cancer.*
14
15 562 *Invest*, 2007, **25**, 397-404.
16
17 563
18
19 564 31. E. Gazi, M. Baker, J. Dwyer, N. P. Lockyer, P. Gardner, J. H. Shanks, R. S. Reeve, C. A. Hart, N. W.
20
21 565 Clarke and M. D. Brown, *Eur. J. Urol*, 2006, **50**, 750-761.
22
23 566
24
25 567 32. P. Lasch, W. Haensch, D. Naumann and M. Diem, *Biochim. Biophys. Acta*, 2004, **1688**, 76-186.
26
27 568
28
29 569 33. H. Martens, J. P. Nielsen and S. B. Engelsen, *Anal. Chem*, 2003, **75**, 394-404.
30
31 570
32
33 571 34. A. Kohler, C. Kirschner and A. Oust, H. Martens, *Appl. Spectrosc*, 2005, **59**, 707-716.
34
35 572
36
37 573 35. D. Sebiskveradze, V. Vrabie, C. Gobinet, A. Durlach, P. Bernard, E. Ly, M. Manfait, P. Jeannesson
38
39 574 and O. Piot, *Lab. Invest*, 2011, **91**, 799-811.
40
41 575
42
43 576 36. P. Bassan, A. Sachdeva, A. Kohler, C. Hughes, A. Henderson, J. Boyle, J. H. Shanks, M. Brown, N.
44
45 577 W. Clarke and P. Gardner. *Analyst*, 2012, **137**, 1370-1377.
46
47 578
48
49 579 37. D. C. Fernandez, R. Bhargava, S. M. Hewitt and I. W. Levin, *Nat. Biotechnol*, 2005, **23**, 469-474.
50
51 580
52
53 581 38. L. M. Wang, D. Kevans, H. Mulcahy, J.O. Sullivan, D. Fennelly, J. Hyland, D. Donoghue and K.
54
55 582 Sheahan, *Am. J. Surg. Pathol*, 2009, **33**, 134-141.
56
57 583
58
59
60

22

- 1
2
3 584 39. H. Kanazawa, H. Mitomi, Y. Nishiyama, I. Kishimoto, N. Fukui, T. Nakamura and M. Watanabe,
4
5 585 *Colorectal. Dis*, 2008, **10**, 41-47.
6
7 586
8
9 587 40. F. Prall, *Histopathology*, 2007, **50**, 151-162.
10
11 588
12
13 589 41. H. Gabbert, *Cancer. Metastasis. Rev*, 1985, **4**, 293-309.
14
15 590
16
17 591 42. R. K. Sahu, S. Argov, S. Walfisch, E. Bogomolny, R. Moreh and S. Mordechai, *Analyst*, 2010, **135**,
18
19 592 538-544.
20
21 593
22
23 594 43. V. K. Katukuri, J. Hargrove, S. J. Miller, K. Rahal, J. Y. Kao, R. Wolters, E. M. Zimmermann and T. D.
24
25 595 Wang, *Biomed. Opt. Express*, 2010, **1**, 1014-1025.
26
27 596
28
29 597 44. B. Rigas, S. Morgello, I. S. Goldman and P. T. Wong, *Proc. Natl. Acad. Sci. USA*, 1990, **87**, 8140-
30
31 598 8144.
32
33 599 45. S. Argov, J. Ramesh, A. Salman, I. Sinelnikov, J. Goldstein, H. Guterman and S. J. Mordechai, *J.*
34
35 600 *Biomed. Opt*, 2002, **7**, 1-7.
36
37 601
38
39 602 46. P. Wong and H. M. Yazdi, *Appl. Spectrosc*, 1993, **44**, 1830-1836.
40
41 603
42
43 604 47. S. Argov, R. K. Sahu, E. Bernshtain, A. Salman, G. Shohat, U. Zelig and S. Mordechai, *Biopolymers*,
44
45 605 2004, **75**, 384-392.
46
47 606
48
49 607 48. S. L. Patrick, T. T. Wong and H. M. Yazdi, *Appl. Spectrosc*, 1993, **47**, 1830-1836.
50
51 608
52
53 609 49. L. Chen, H. Y. N. Holman, H. Zhao, H. A. Bechtel, M. C. Martin, C. Wu and S. Chu, *Anal. Chem*, 2012, **84**,
54
55 610 4118-4125.
56
57
58
59
60

23

611 **Figures:**

612

613

614 **Figure 1:** K-means clustering and digital staining of FTIR spectral images with random pseudo-colors.

615 Left panel: HPS stained colon tissues (Supplementary Table S1, sample # 1D and 12C); Middle panel: K-

616 means clustering and digital staining of FTIR spectral images with random pseudo-colors; Right panel:

617 Dendrograms corresponding to the respective cluster images. A is a non-tumoral colonic tissue

618 partitioned using 8 clusters representing the major normal colonic tissue features. The cluster

619 representation is as follows: Cluster 1 - lamina propria, cluster 2 - mucous, clusters 4, 5 and 7 -

620 submucosa, cluster 6 - crypt (outer part-OP), cluster 8 - crypt (inner part IP) and cluster 3 - undefined

621 tissue. B is a moderately differentiated colonic adenocarcinoma partitioned using 11 clusters. The

622 important histological classes are cluster 1 - tumor, clusters 6, 7, and 11 - tumor-associated stroma.

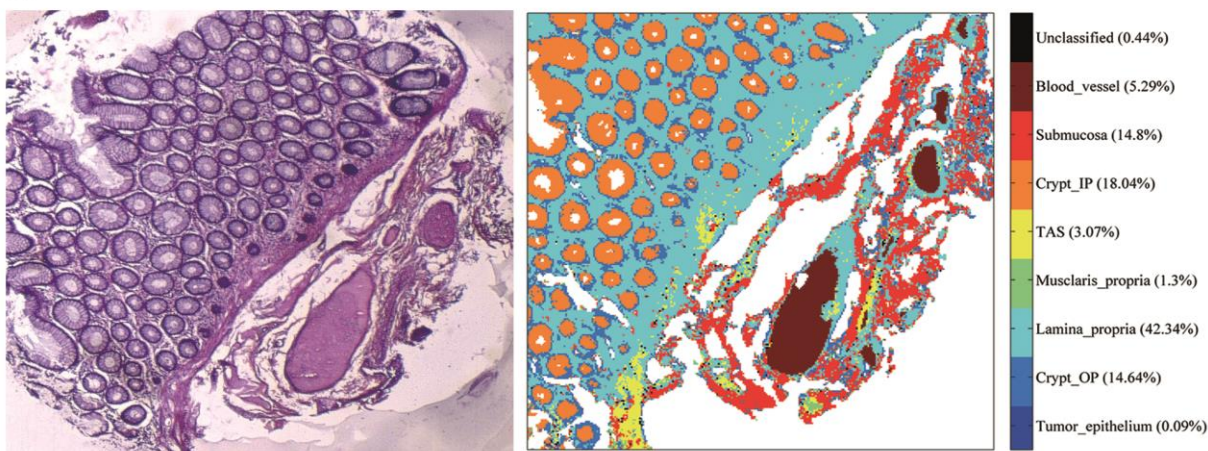
623 Remaining clusters were attributed to the fibrous stroma. The HPS images are at 5X magnification.

624

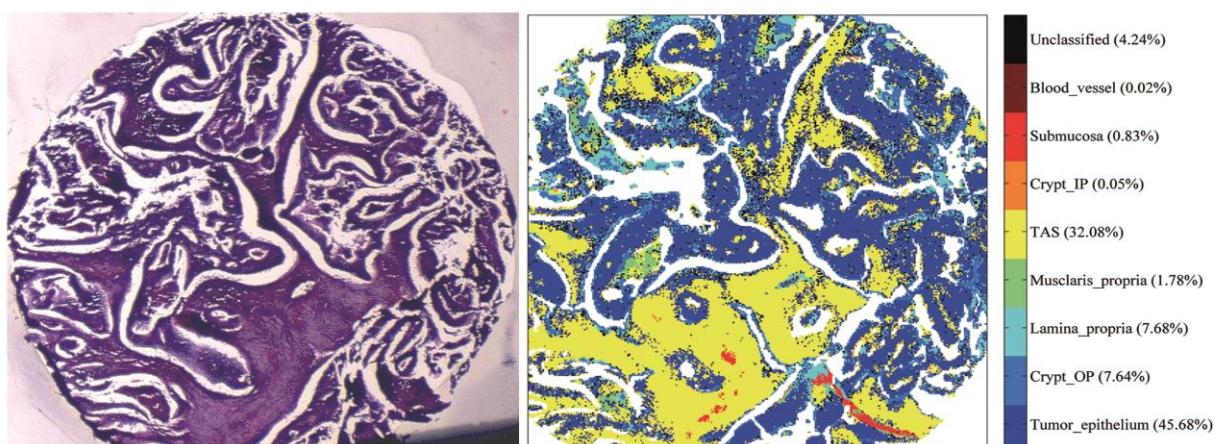
625

24

A



B



626

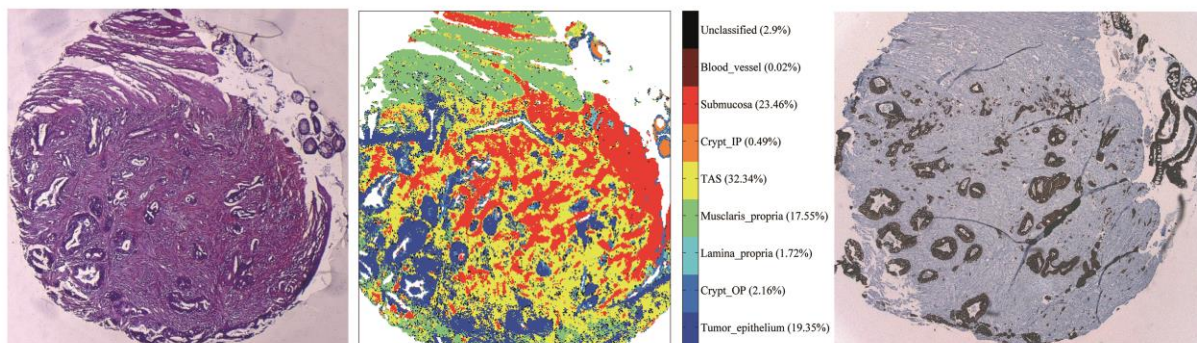
627

628 **Figure 2:** Performance of the prediction model: Identification of unknown colonic tissues by spectral
 629 histopathology. Left panel: HPS stained colon tissues (Supplementary Table S1, sample # 14D and 7C);
 630 Right panel: Infrared spectral predicted images. A is a non-tumoral colonic tissue section in which all the
 631 important normal colonic histological features are well-identified by the model. The important histological
 632 classes such as normal epithelium (crypt-IP and crypt-OP), connective tissue, blood vessels, etc are
 633 represented by a specific color-code. B is a moderately differentiated colon adenocarcinoma in which the
 634 tumor epithelium is together with its associated stroma are represented by the specific color-code. Note
 635 that there is a complete absence of normal epithelium. The HPS images are at 5X magnification.

636

25

637



638

639

640 **Figure 3:** Identification of tumor budding in an unknown colonic tissue. Left panel: HPS stained colonic
641 tissue (Supplementary Table S1, sample # 9B); Middle panel: Infrared spectral predicted image; Right
642 panel: KL 1 immuno-stained image. The sample is a moderately differentiated colon adenocarcinoma in
643 which the cancerous glands are identified along with the tumor-associated stroma. Small isolated tumor
644 clusters representing tumor-budding are identified branching out into the stroma. The tumor-stromal
645 boundary is also well-identified and clearly demarcated from the normal connective tissue (muscularis
646 propria). In the same sample, few normal colonic glands are seen in the top-right position identified by
647 presence of normal epithelium. The HPS and the IHC images are at 5X magnification.

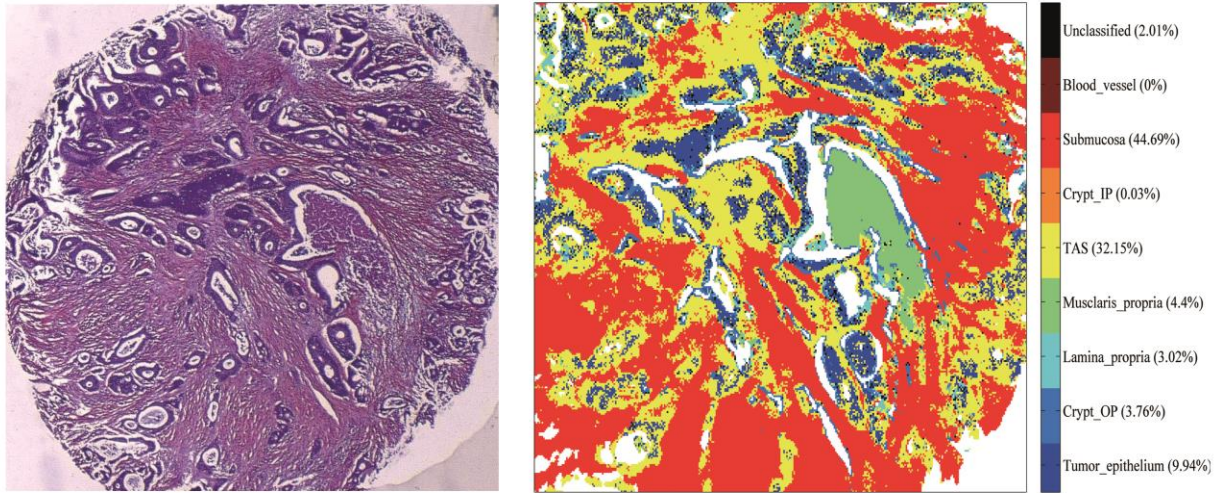
648

649

650

651

26



652

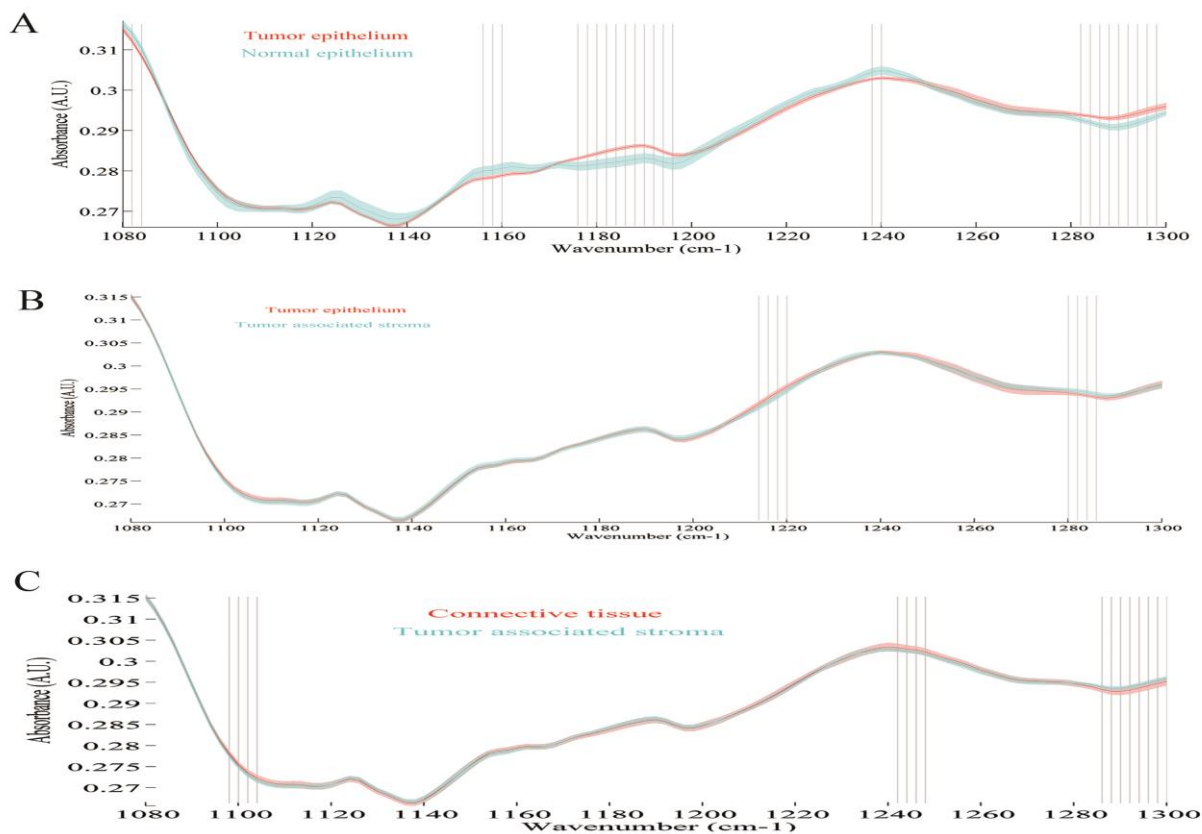
653

654

655 **Figure 4:** Tumor stroma geographical proximity. The sample is a moderately differentiated colonic
656 adenocarcinoma with its associated stroma (Supplementary Table S1, sample # 11B). Along with the
657 highly-correlated prediction, the nature of the connective tissue into which the tumor has infiltrated is also
658 identified. The HPS image is at 5X magnification.

659

27



660

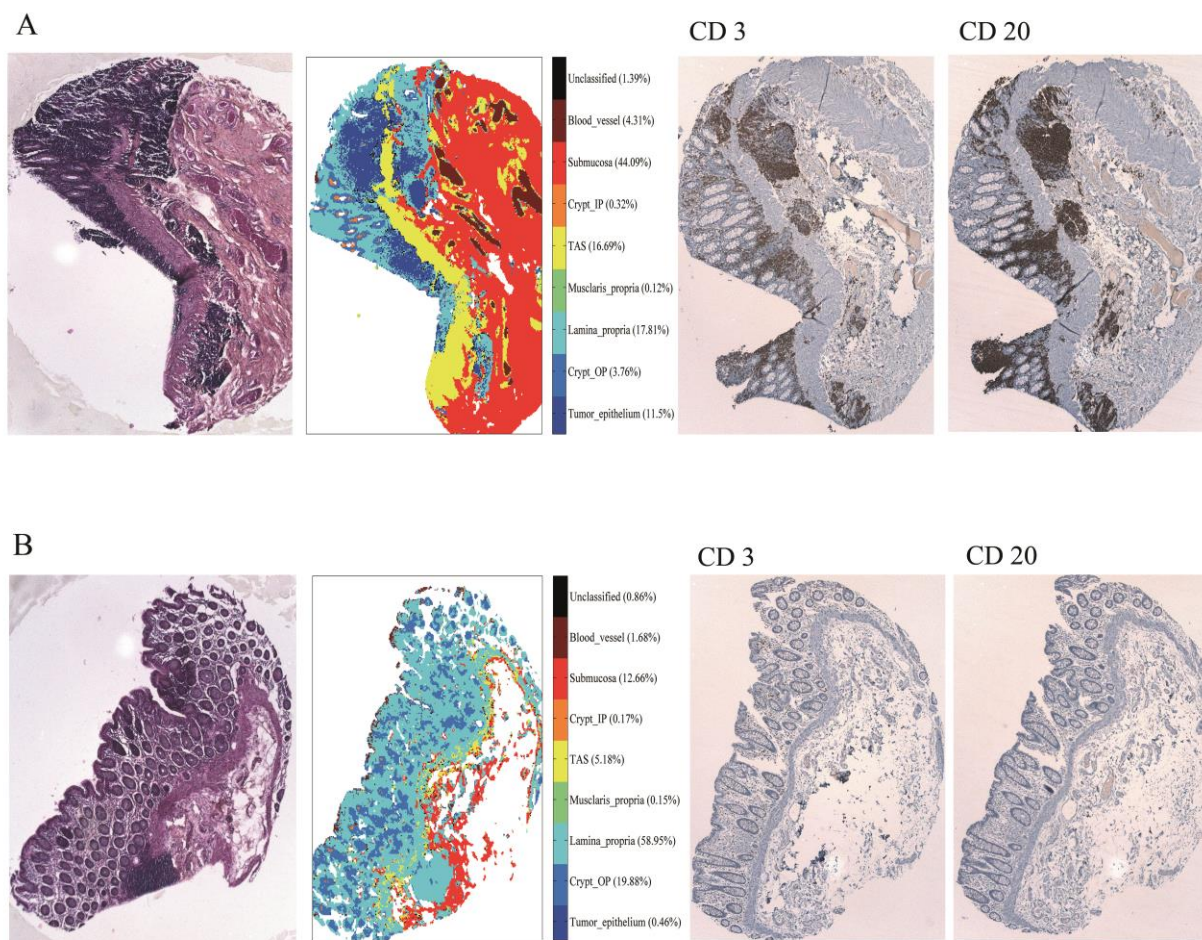
661

662 **Figure 5:** Most discriminant infrared spectral vibrations identified by the Mann-Whitney U test. The test
663 was performed for A: Tumor epithelium versus normal epithelium ($p < 0.005$), B: Tumor epithelium versus
664 tumor associated stroma ($p < 0.01$), and C: Connective tissue versus tumor associated stroma ($p < 0.1$). For
665 each class in the figure, the mean spectrum (+/-) the standard deviation is represented.

666

667

28



668

669 **Figure 6:** Influence of tissue inflammation on the prediction model. Left to right: 1. HPS stained colon

670 tissues (Supplementary Table S1, sample # 32 and 31); 2. Infrared spectral predicted images; 3.

671 Immuno-stained images for CD3 marker and; 4. Immuno-stained images for CD20 marker. A is a non-

672 tumoral colonic tissue with typical normal glands. The mucosa is partially populated by lymphoid follicle as

673 seen in the HPS image. The prediction model identified these regions in the mucosa as tumor. Immuno-

674 staining for CD3 and CD 20 markers revealed that the tumor class in the predicted images actually

675 corresponded to inflammatory signatures. B is another non-tumoral tissue with insignificant tumor pixels in

676 the predicted image. Immuno-staining is negative for CD 3 and CD 20 indicating absence of inflammatory

677 signature. The HPS and the IHC images are at 5X magnification.

678

29

679 **Tables:**

680

681 **Table 1:** The confusion matrix.

682 The confusion matrix representing the sensitivity of the infrared spectral imaging based prediction model,
 683 developed using 8 classes, to the gold standard histopathological attribution, in the spectral range of 1080
 684 cm^{-1} to 1300 cm^{-1} . The table shows an average sensitivity of 89.49 %.

685

686

Predicted class (infrared imaging)

	'Tumor_epithelium'	'Crypt_OP'	'Lamina_propria'	'Musclaris_propria'	'TAS'	'Crypt_IP'	'Submucosa'	'Blood_vessel'
'Tumor_epithelium'	96,45	0	1,32	0,07	2,03	0	0	0,1
'Crypt_OP'	0,1	88	6,16	0,28	0	4,42	0,5	0,22
'Lamina_propria'	2	1,27	83	0,14	11,34	0,09	1,8	0,45
'Musclaris_propria'	0,1	0,04	0	98,22	1,1	0,04	0	0,04
'TAS'	16,33	0	1	0,08	81,31	0,02	1,24	0
'Crypt_IP'	0,04	6	0,04	0,04	0	93,7	0,16	0,04
'Submucosa'	0,1	0	7,5	0,05	14,42	0	77,54	0,35
'Blood_vessel'	0	0	0	0	2,3	0	0	97,7

No. of spectra used in the model	35083	3567	14106	4514	16409	8377	3964	782
----------------------------------	-------	------	-------	------	-------	------	------	-----

Total = 86802

687

688

689

690

691

692

693

694