

1 Prediction of malignant transformation in oral
2 epithelial dysplasia using infrared absorbance
3 spectra
4

5 Barnaby G. Ellis¹, Conor A. Whitley¹, Asterios Triantafyllou², Philip J. Gunning³,
6 Caroline I. Smith¹, Steve D. Barrett¹, Peter Gardner⁴, Richard J. Shaw^{*3,5}, Peter
7 Weightman¹ and Janet M. Risk³
8

9 ¹ Department of Physics, University of Liverpool, UK.

10 ² Department of Pathology, Liverpool Clinical Laboratories, University of
11 Liverpool, UK.

12 ³ Department of Molecular and Clinical Cancer Medicine, Institute of Systems,
13 Molecular and Integrative Biology, University of Liverpool, UK.

14 ⁴ Manchester Institute of Biotechnology, University of Manchester, UK.

15 ⁵ Regional Maxillofacial Unit, Liverpool University Hospitals NHS Foundation
16 Trust, Liverpool, UK.
17

18 * Corresponding author

19 E-mail: rjshaw@liverpool.ac.uk (RJS)

20

21 **Abstract**

22 Oral epithelial dysplasia (OED) is a histopathologically-defined, potentially
23 premalignant condition of the oral cavity. The rate of transformation to frank
24 carcinoma is relatively low (12% within 2 years) and prediction based on
25 histopathological grade is unreliable, leading to both over- and under-treatment.
26 Alternative approaches include infrared (IR) spectroscopy, which is able to classify
27 cancerous and non-cancerous tissue in a number of cancers, including oral. The aim
28 of this study was to explore the capability of FTIR (Fourier-transform IR)
29 microscopy and machine learning as a means of predicting malignant transformation
30 of OED. Supervised, retrospective analysis of longitudinally-collected OED biopsy
31 samples from 17 patients with high risk OED lesions: 10 lesions transformed and 7
32 did not over a follow-up period of more than 3 years. FTIR spectra were collected
33 from routine, unstained histopathological sections and machine learning used to
34 predict malignant transformation, irrespective of OED classification. PCA-LDA
35 (principal component analysis followed by linear discriminant analysis) provided
36 evidence that the subsequent transforming status of these 17 lesions could be
37 predicted from FTIR data with a sensitivity of $79 \pm 5\%$ and a specificity of $76 \pm 5\%$.
38 Six key wavenumbers were identified as most important in this classification.
39 Although this pilot study used a small cohort, the strict inclusion criteria and
40 classification based on known outcome, rather than OED grade, make this a novel
41 study in the field of FTIR in oral cancer and support the clinical potential of this
42 technology in the surveillance of OED.

43

44 Introduction

45 Oral squamous cell carcinoma (OSCC) has a worldwide incidence rate of
46 over 370 000 [1] and a 5-year survival rate that remains less than 60% [2]. It is often
47 preceded by a spectrum of clinical changes, collectively termed potentially pre-
48 malignant oral epithelial lesions or oral potentially malignant disorders (OPMDs)
49 [3]. OPMDs include white patches (leukoplakia), red patches (erythroplakia) and
50 mixed color patches (erythroleukoplakia) of the oral mucosa and are usually
51 investigated by biopsy and histopathological examination. The latter typically shows
52 morphological changes of the surface squamous epithelium, which are conveniently
53 described as oral epithelial dysplasia (OED), that include variable cellular atypia,
54 proliferative activity and loss of normal patterns of differentiation. The World Health
55 Organization (WHO) defines OED as “altered epithelium with an increased
56 likelihood for progression to squamous cell carcinoma” [4]. The malignant potential
57 (progression to invasive tumor) of OPMDs range from as low as 0.13% in some
58 leukoplakias [5] to >50 % in some erythroplakias [6]; a meta-analysis of OED data
59 indicates a malignant transformation rate of 12% within 2 years, increasing to 22%
60 within 5 years [7].

61 The histopathological grading of OED as mild, moderate or severe is widely
62 popular and time-honored, based on various architectural and cytological changes,
63 and is endorsed and formalized by WHO [8]. Other grading systems, including a
64 binary “high/low risk” scheme [9], have been proposed, but the standard across the
65 field remains the three-tiered scheme. Numerous studies [10-12] have shown a
66 significant relationship between the histopathological grade and risk of malignant
67 transformation, however there are a similar number of conflicting reports which
68 suggest a much less direct relationship, highlighting other risk factors [13-15].
69 Although biopsy and histopathological assessment of OPMDs forms the basis of
70 clinical management, the grading of OED is influenced by inter- and intra-observer
71 variations [9], reflecting the subjectivity of the process, and improvements are
72 required.

73 Given its clinical significance, OED has been investigated by a wide
74 spectrum of predominantly biology-based methodologies, but none of the proposed
75 biomarkers for predicting risk are in routine clinical use [16]. Less attention has been
76 paid to the application of alternative methodologies utilizing the chemical or
77 physical properties of the cells. Among the alternative methodologies, those based
78 on vibrational spectroscopy have been increasingly introduced to biomedical
79 research [17, 18]. Fourier transform infrared (FTIR) spectroscopy utilizes infrared
80 (IR) light over a broad spectral range to assess the overall chemical profile of a
81 sample. Molecules which vibrate at frequencies corresponding to the wavelengths
82 applied will absorb the radiation at those wavelengths, resulting in an absorption
83 spectrum characteristic of the chemical moieties present. FTIR micro-spectroscopy
84 (FTIR-MS) combines IR spectroscopy with precise spatial information enabling the
85 rapid acquisition of hyperspectral images directly related to the location and
86 distribution of chemical components, for example in tissue samples. Hyperspectral

87 data acquired using FTIR-MS is highly dimensional as each raw spectrum, obtained
88 from a region approximately $5\ \mu\text{m} \times 5\ \mu\text{m}$ in size, contains at least 10^3 absorption
89 variables. Subtle differences between tissue areas are concealed by dominant
90 common features in chemical composition, necessitating the use of sophisticated
91 numerical approaches to extract useful information [19]. Common modelling
92 methods whose aim is to reduce the complexity of the dataset include principal
93 component analysis (PCA) [20] and linear discriminant analysis (LDA) [21, 22].

94 FTIR-MS has been utilized in biomedical research, with a particular focus
95 on its application to the investigation of cancerous tissues (reviewed in ref. [23]).
96 Our own recent data suggests that this methodology is applicable to OSCC [24], and
97 other studies have successfully associated vibrational spectroscopy data with the
98 contemporaneous histopathological classification of potentially malignant oral
99 lesions [25, 26]. The present investigation takes a different approach, using a
100 supervised, retrospective analysis of tissue samples from high risk OED lesions from
101 patients with prolonged, longitudinal clinical follow-up and known outcome
102 (transformation or no transformation) to explore the capability of FTIR-MS and
103 machine learning as a means of predicting malignant transformation of OED.

104

105 **Methods**

106 Seventeen patients with biopsy-proven OED were included in this study
107 (Table 1). All patients were part of a larger cohort for whom the clinical determinants
108 of transformation have been described [15] and had given written informed consent
109 to a UK NHS Research Ethics Committee approved study that was run in compliance
110 with the Helsinki Declaration (Liverpool Central REC ref: EC 47.01). Patient
111 selection was limited by inclusion of only lesions with a histopathological diagnosis
112 of moderate or severe grade OED, absence of previous OSCC, at least 42 months
113 follow-up from the time of biopsy, and the availability of relevant archival formalin-
114 fixed paraffin-embedded (FFPE) tissue. Although inclusion was partly dependent on
115 a number of non-clinical factors such as availability of samples, the group of patients
116 used in this present study remained representative of the total cohort [15]. Thus, there
117 was a higher female:male ratio in the transforming group, which also had a
118 preponderance of lateral tongue lesions, and there were proportionally more smokers
119 in the non-transforming group (Table 1). It should be noted that there was no
120 significant difference between the distribution of severe and moderate grade OED
121 lesions in the two groups.

122

123 **Table 1. Patient and Sample cohort characteristics.**

Patient group	Identifier	Age at biopsy	Gender	Site	Number of clinical sites	Clinical Presentation	Lesion size (mm ²) at presentation	Histology Grade of ROI ^a	Time before transformation (months)	Time cancer free (months)	Lifestyle	
											tobacco	alcohol
Transforming n=10	12089	58	M	Ventral Tongue	single	Erythroleukoplakia	101-500	severe	2		y (2-20py) ^b	y
	12201	52	M	Buccal	single	Leukoplakia	>500	severe	4		y (>20py)	y
	12260	74	F	Floor of Mouth	single	Leukoplakia	≤100	moderate-severe	5		y (>20py)	y
	12127	85	M	Lateral Tongue	single	Leukoplakia	100-500	moderate	7		n	n
	12257	49	F	Lateral Tongue	single	Leukoplakia	100-500	moderate	12		n	y
	12248	69	F	Lateral Tongue	multiple	Erythroleukoplakia	100-500	severe	14		n	n
	12263	45	F	Lateral Tongue	single	Leukoplakia	100-500	moderate-severe	18		y (5-20py)	y
	12104	70	F	Lateral Tongue	single	Erythroleukoplakia	100-500	moderate	26		n	n
	12195	45	F	Ventral Tongue	multiple	Erythroleukoplakia	>500	moderate	33		y (5-20py)	y
	12219	68	F	Ventral Tongue	single	Leukoplakia	≤100	moderate-severe	43		y (5-10py)	n
Non-Transforming n=7	12181	49	F	Soft Palate	single	Leukoplakia	>500	severe		158	n	y
	12330	47	F	Soft Palate	multiple	Erythroleukoplakia	>500	moderate		108	y (5-20py)	y
	12098	78	M	Ventral Tongue	single	Leukoplakia	>500	moderate		106	y (>20py)	y
	12332	71	F	Mandibular alveolus	multiple	Leukoplakia	≤100	moderate		91	n	y
	12329	59	M	Floor of Mouth	single	Erythroleukoplakia	≤100	moderate		75	y (>20py)	y
	12162	61	F	Floor of Mouth	single	Erythroleukoplakia	≤100	severe		67	y (>20py)	y
	12141	47	F	Floor of Mouth	multiple	Erythroleukoplakia	≤100	severe		43	y (>20py)	y

124 ^aROI = region of interest: the target for FTIR spectroscopy;

125 ^bpy=pack years (= packs per day multiplied by years smoked)

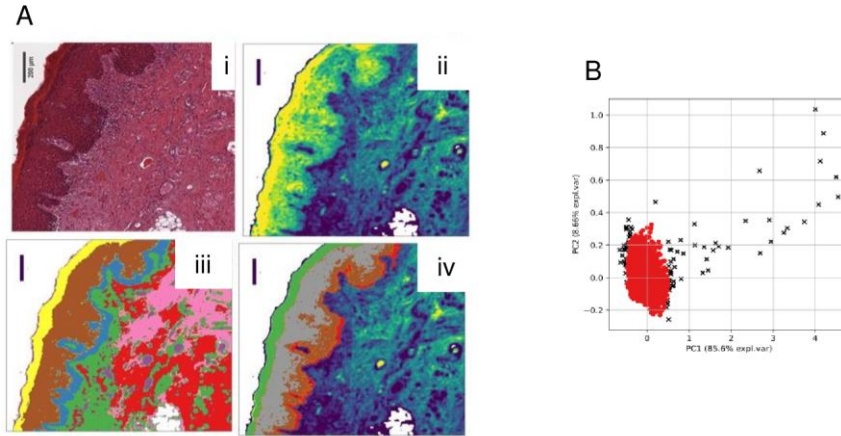
126 A single archival FFPE tissue block containing incisional biopsy material
127 was obtained from each of 10 patients at the closest timepoint to transformation
128 (range 2-43 months prior to transformation) (T lesions; Table 1). A single archival
129 FFPE block containing incisional biopsy material with more than 43 months
130 transformation-free follow-up from the date of biopsy (range 43-108 months) was
131 obtained from each of 7 patients (NT lesions; Table 1). None of these NT lesions
132 had been excised during the follow-up period.

133 From each of the 17 FFPE blocks, four adjacent 5 μm tissue sections were
134 obtained, reserving the first and last for routine deparaffinization, hematoxylin and
135 eosin (H&E) staining and histopathological re-examination. The intervening two
136 paraffinized, unstained sections were mounted on separate 20 mm diameter calcium
137 fluoride (CaF_2) disks for IR imaging experiments. An area of dysplasia
138 corresponding to the histopathologically most extreme OED present was identified
139 in each H&E-stained section and marked as the target for FTIR spectroscopy. This
140 area was termed the region of interest (ROI).

141 FTIR imaging data were acquired using a Varian 620 microscope coupled to
142 an Agilent Cary 670 spectrometer (Agilent, Stockport, UK) enclosed within a
143 purging chamber to eliminate water vapor and carbon dioxide contributions. The
144 instrument was configured to collect mid-IR transmission data between 900 and
145 3800 cm^{-1} with a spectral resolution of 4 cm^{-1} and pixel size of $5.5\text{ }\mu\text{m}$, allowing the
146 simultaneous acquisition of 128×128 spectra over a field of view of approximately
147 0.5 mm^2 . CaF_2 disks with mounted sections were loaded onto a 3D-printed slide
148 holder capable of containing three disks. Disks were imaged two at a time, with the
149 third position reserved for a clean, blank disk to allow for microscope calibration
150 and spectral background correction. The semiconductor detector (Mercury-
151 Cadmium-Telluride) in the FTIR microscope was cooled with liquid nitrogen to
152 78 K in order to reduce thermal noise in the data. Using the average of 128 scans of
153 the blank disk for background correction, a hyperspectral IR image was obtained for
154 each ROI by averaging 64 scans of the identified dysplastic area. A built-in mosaic
155 function was utilized for cases where the extent of the surface oral epithelium in the
156 tissue sections could not be visualized in one field of view of the microscope and
157 enabled the acquisition of larger, composite images.

158 Subsequently, the FTIR images were cross-referenced with scanned images
159 of the corresponding H&E sections to confirm the location and extent of dysplasia
160 within the ROI (Fig 1A). To annotate IR spectral data originating in dysplastic
161 epithelium, each hyperspectral image was subjected to a two-tiered, k-means cluster
162 analysis (Fig 1A). Initially, the epithelium was identified by a specialist, head and
163 neck histopathologist (AT). The first clustering step was then used to identify this
164 structure based on its IR hyperspectral profile. The data from this first clustering step
165 was then processed during the second clustering in order to identify regions within
166 the epithelium based on their IR hyperspectra. Histologically, dysplasia was often
167 centered on the basal and parabasal layers and, in the case of more severe dysplasia,
168 in the upper prickle-cell layer. IR data from areas with both a histological assessment

169 of dysplasia and where k-means clustering identified relative chemical homogeneity
170 were selected for modelling.
171



172

173 **Fig 1. Identification of IR data to be used in classification.** (A) Example of two-tiered k-
174 means cluster analysis. (i): H&E image; (ii): corresponding FTIR hyperspectral image; (iii):
175 the first tier of k-means cluster analysis identifies the surface epithelium as 3 separate,
176 spectrally similar regions (identified as yellow, brown and blue colored layers by
177 histopathological comparison with (i)). Histologically, the blue colored cluster broadly
178 corresponds to the basal layer; the brown cluster to parabasal and prickle-cell (spinous)
179 layers; and the yellow to the keratinized layers; (iv) the second tier of k-means cluster
180 analysis subdivides the epithelium into four clusters of spectrally similar regions (green,
181 grey, brown & orange). This second 2-tier clustering appears to separately identify the
182 parabasal (brown) and spinous (grey) layers. Histopathology plus PCA clustering of FTIR
183 data selects the brown and red clusters for use in modelling. Scale bar = 200 μ m. (B)
184 Illustration of the quality control process. Spectra identified as lying outside the 95%
185 confidence interval by the Hotelling's T-squared test (black crosses) were removed from
186 dataset. Data in this figure were obtained from the same tissue section as in part (A).

187

188 Spectra that originated from dysplastic material in each IR hyperspectral
189 image were then subject to an initial quality check to discard anomalous spectra.
190 This involved using PCA as a tool to decompose the spectra from each image into
191 five principal components, and then employing Hotelling's T^2 summary statistic to
192 determine which spectra lie furthest away from the origin, discarding any which lie
193 outside the 95% confidence interval [27]. The remaining spectra were retained for
194 modelling. Although the FTIR and H&E images had been cross-referenced in order
195 to locate the imaged dysplastic epithelial layers, the spectral data were grouped into
196 two categories based on the clinical outcome of the lesion from which they were
197 taken regardless of OED severity: T lesions underwent malignant transformation and
198 NT lesions did not undergo transformation.

199 We have developed an objective optimization framework, PipeOpt (paper in
200 preparation), that aims to maximize the efficiency of a classifier by optimizing the
201 parameters of each pre-processing step of the IR data, determining their ideal
202 sequence and identifying the best performing classification method(s) using
203 Bayesian optimization (Table 2). This process probabilistically converges on the best

204 hyperparameters for each unique series of pre-processing and classification steps
205 (defined as a pipeline) by iteratively updating the associated Matthew's correlation
206 coefficient [28, 29]: a statistic scaled between -1 and 1 which includes consideration
207 of both the sensitivity and specificity of the resulting classification. This allows for
208 faster convergence to the ideal compared to an approach that samples every possible
209 combination of hyperparameters. To assess the performance of each pipeline, 70
210 training sets were created by using a leave-one-pair-out cross validation (LOPOCV)
211 method: pairs of samples (comprising 1 T and 1 NT sample in every possible
212 combination) was set aside as the test set, while data from the remaining 6 NT and
213 9 T samples was used for training. Equal numbers of spectra (n=500) were used from
214 each sample to avoid sample-related bias that might influence the optimization and
215 the optimal pipeline was defined as that with the highest mean Matthew's correlation
216 coefficient.

217 Following determination of the optimal pipeline, the same sequence of pre-
218 processing and classifier steps were used to analyze all of the available data (range
219 3891-5437 spectra) from the 17 samples using the same LOPOCV routine as before
220 to create training and test sets.

221 PCA-LDA is a feature extraction and classifier hybrid that uses a series of
222 linear transformations to decompose the data from absorption variables to a single
223 variable called a linear discriminant (LD), the value of which is called the LD score.
224 The LD score is dependent upon the PCA and LDA loading vectors, which are a
225 measure of the relative weight that each wavenumber in the spectrum contributes to
226 that part of the analysis. The LD scores for each datapoint in every lesion were
227 plotted against their frequency of occurrence in relation to their known
228 transformative capacity and were also used to identify key wavenumbers in the
229 discrimination of T from NT datapoints.

230 At each iteration of train/test, the predicted outcome (transformation or no
231 transformation) at every datapoint for the test sample was determined and the lesion
232 as a whole classified as T or NT based on a simple majority.

233 The FTIR data was converted from native data format using the ChiToolbox
234 package MATLAB [30]. All data transformation and statistical analyses were
235 performed using either in-house developed packages (PipeOpt) or third-party
236 packages implemented in Python v3.9.

237
238

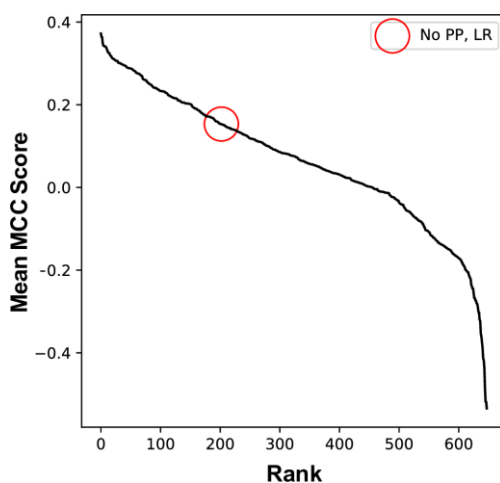
239 **Table 2. Pre-processing steps tested in the PipeOpt objective optimization**
 240 **framework.**

Step	Method	Hyperparameter	Options
Smoothing	Savitzky-Golay (SG) smoothing	Window size	5,7,9,11,13,15,17,19,21
	PCA	Explained Variance	80-95%
	none		
Baseline Correction	Rubberband	N/A	Y/N
	SG Differentiation	Window size (if no smoothing)	5,7,9,11,13,15,17,19,21
		Polynomial order	2,3
		Differentiation order	1,2
none			
Paraffin Correction	Removal of the spectral region dominated by paraffin wax (1340 – 1490 cm ⁻¹)		
Normalisation	Vector	N/A	Y/N
	Min-max	N/A	Y/N
	Amide I	N/A	Y/N
	none		
Scaling	Standard scaling	N/A	Y/N
	Min-max	N/A	Y/N
	none		
Feature Extraction	PCA	Explained Variance	90-98%
	none		
Classifier	Logistic regression	Regularisation strength	0.001-10
	Linear discriminant analysis	N/A	Y/N
	Random Forest	Max-depth	Y/N
		Minimum samples per split	2,3,4,5
		Minimum samples per leaf	1,2
Bootstrap		Y/N	

241 Legend: Step: typical pre-processing step utilized in analyzing IR data;
 242 Method: methods typically used to perform each pre-processing step;
 243 Hyperparameter: typical parameters associated with each method;
 244 Options: typical options for each parameter. Each step (apart from Paraffin
 245 correction and Classification) also has a bypass option and the steps may be
 246 performed in any order (except for classification).
 247 Bayesian optimization was used to identify hyperparameter options for the resulting
 248 3 x 3 x 1 x 4 x 3 x 2 x 3 = 648 pipelines.

249 Results

250 Each of 648 different data pipelines were tested 70 times, each iteration
251 having a different pair of samples (one T and one NT) removed to create the test/train
252 datasets, and the mean Matthew's correlation coefficient across these 70 iterations
253 was determined (Fig 2). The pipeline with the highest mean correlation coefficient
254 (0.37) was identified from this analysis and is as follows: denoising of the spectra
255 using the Savitzky-Golay smoothing algorithm [31] with a window size of 15 and
256 polynomial order of 2; first order differentiation of the spectra to remove effects such
257 as scattering and background interferences; removal of the spectral region dominated
258 by paraffin wax ($1340 - 1490 \text{ cm}^{-1}$) [23]; normalization so that the sum of the
259 squares of each spectrum is equal to 1 (vector normalization) to account for
260 variations in sample thickness; PCA-LDA classification. In this analysis, PCA was
261 applied to the spectral data to decompose it into the number of principal components
262 that described 90% of the explained variance in the original dataset and LDA was
263 then used to discriminate between the two groups of lesions (T and NT) using the
264 principal components as input.
265



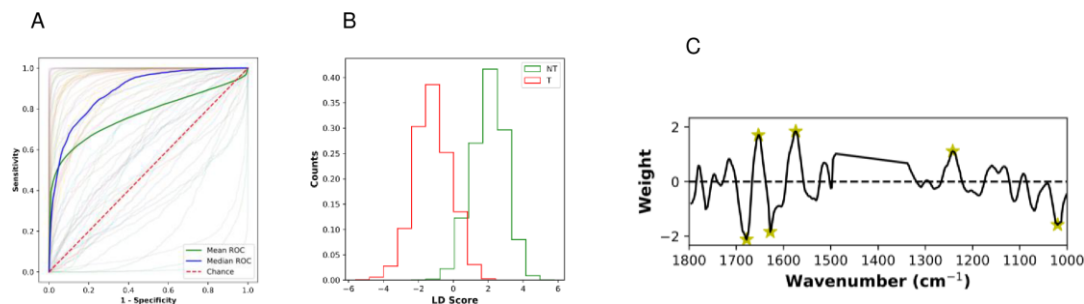
266

267 **Fig 2. Determination of the optimal analysis pipeline for this dataset.** Mean of
268 Matthew's correlation coefficient (MCC) for each of 648 analysis pipelines generated from
269 the dysplasia dataset plotted in descending order. Circle identifies the MCC for the pipeline
270 with no data pre-processing and classification using linear regression (MCC=0.15).
271

271

272 The mean sensitivity of this discriminatory model when applied to the whole
273 dataset at the level of an individual spectrum (i.e. each datapoint in every sample
274 taken as an individual element) was $74 \pm 2.8\%$ and the specificity was $69 \pm 3.2\%$,
275 while the mean and median receiver operator characteristic (ROC) further
276 demonstrated the performance of the model (Fig 3A). Moreover, the PCA-LDA-
277 derived linear discriminant score showed good separation of the two classes (T and
278 NT) (Fig 3B) and the weighting assigned to different wavenumbers during this
279 analysis, allowed the provisional identification of six wavenumbers that provided
280 the most discriminatory power (Fig 3C).

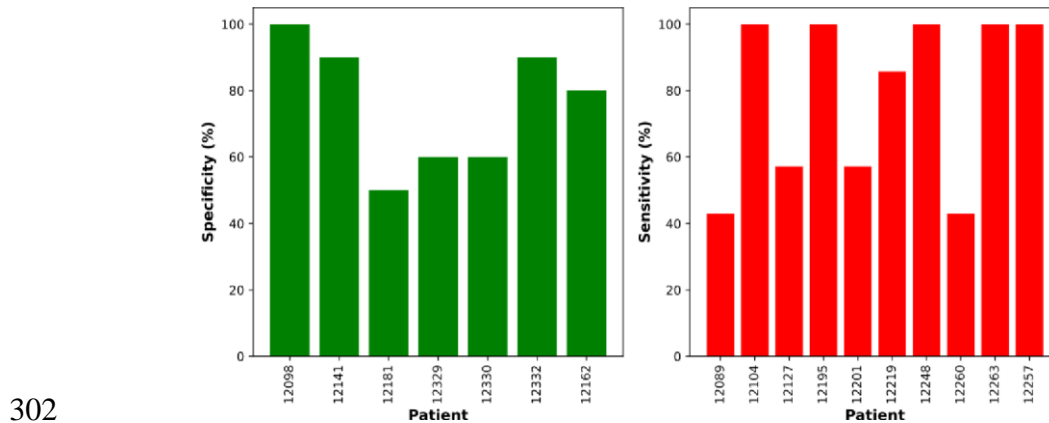
281



282 **Fig 3. Performance of the model taking each datapoint individually.** Means from 70
283 iterations of train/test sets are presented. (A) Mean (green) and median (blue) receiver
284 operator characteristic (ROC) curves. Red dotted line would be achieved by random chance.
285 Pale lines are individual ROC curves for each iteration; (B) Histogram of frequency of
286 occurrence of linear discriminant (LD) scores (counts) plotted for all datapoints from
287 transforming (T: red) and non-transforming (NT: green) lesions, showing separation of the
288 two classes; (C) Plot showing the weighting (a measure of relative importance) assigned to
289 each wavenumber during the PCA-LDA analysis. Features marked with a yellow star show
290 the largest magnitude in weighting: 1678 cm⁻¹, 1653 cm⁻¹, 1628 cm⁻¹, 1574 cm⁻¹, 1242 cm⁻¹
291 and 1020 cm⁻¹.

292

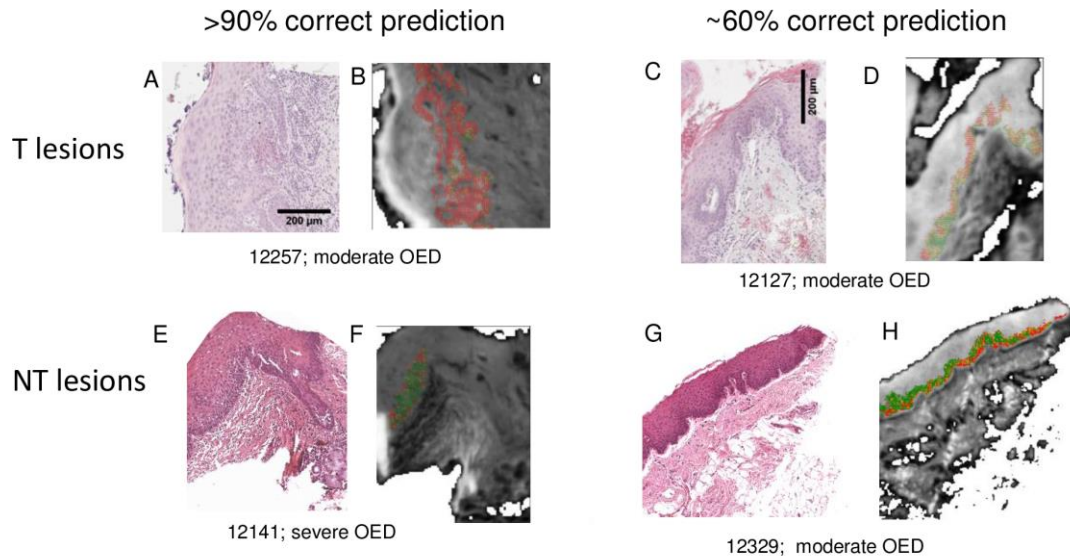
293 However, for clinical utility, the transformative capacity of the whole lesion
294 is more relevant than that for each individual datapoint. Therefore, in each iteration
295 of train/test, the predicted outcome (transformation or no transformation) at every
296 datapoint for the test sample was used to define the lesion as T or NT based on a
297 simple majority (i.e. $\geq 50\%$ of datapoints). The sensitivity per lesion was $79 \pm 4.9\%$
298 and the specificity was $76 \pm 5.1\%$. However, the prognosis of some lesions was
299 better predicted than others (Fig 4), with 2 T lesions and 1 NT lesion being
300 incorrectly predicted in $\geq 50\%$ of the test/train iterations in which they were the test
301 sample.



302

303 **Fig 4. Performance of the model at the lesion level.** Means from 70 iterations of train/test
304 sets are presented. Frequency at which each whole lesion was correctly predicted when it
305 appeared in the test set (based on LOPOCV, each NT lesion is present in 10 test pairs and
306 each T lesion in 7 test pairs). Data is shown for each individual NT (left) and T (right) lesion
307 and plotted as specificity (true negatives) and sensitivity (true positives). Patient numbers
308 are research sample IDs.

309 To better visualize this observation, the probability of transformation for
310 each datapoint from 4 lesions was color coded and mapped back onto a
311 representation of the whole section (Fig 5). Lesions that were most accurately
312 predicted showed homogeneous areas of correctly labelled datapoints with only a
313 few incorrectly labelled points (Fig 5B and F). Conversely, lesions that were less
314 accurately predicted demonstrated some areas that were predicted to transform and
315 some that were predicted to not transform (Fig 5D and H). The OED grade of the
316 lesions did not appear to correlate with the success or otherwise of the prediction of
317 transformation.



318

319 **Fig 5. Representative images demonstrating identification of lesional areas that are**
320 **mis-labelled.** Every datapoint for each lesion was color coded to represent the probability
321 of transformation (T: red; NT: green; see color bar) and mapped back onto a representation
322 of the whole section for two T (top) and two NT (bottom) lesions. Lesions predicted with
323 high (left) and lower (right) accuracy are shown. (A), (C), (E) and (G): H&E images; (B),
324 (D), (F) and (H): corresponding maps of predicted datapoints on IR images.
325

326 Discussion

327 We have applied machine learning to infrared data collected by FTIR-MS
328 from a number of high risk OED lesions with known transforming potential. This is
329 in contrast to most other studies analyzing FTIR data that has been collected from
330 oral premalignant lesions, which commonly correlate IR data with OED stage rather
331 than outcome. The process correctly predicted the capacity to transform with a
332 sensitivity of $79 \pm 4.9\%$ and a specificity of $76 \pm 5.1\%$, which is better than when
333 OED grade alone is used. A histological grade of severe OED or carcinoma-in-situ
334 has been observed to have a significantly increased malignant transformation rate
335 compared with mild or moderate OED ($P < 0.008$) [7], but this grading is still only
336 predictive for 24-40% of such lesions [7, 32]. It is of note that, in the limited cohort

337 used in the present study, equivalent numbers of severe and moderate OED were
338 present in the T and NT groups, and OED grade did not correlate with the ability of
339 the classifier to accurately predict transformation potential.

340 Two lesions known to transform were predicted to be non-transforming in
341 the current analysis. Based on our knowledge of oral cancer development, it may be
342 hypothesized that the lesions are most probably heterogeneous with areas possessing
343 transforming potential and areas without this potential. Thus, the biopsy may not
344 have been representative of the region that underwent subsequent transformation.
345 Similarly, given the relatively large size of the region used for IR imaging compared
346 to the size of individual cells, the IR imaged area will contain areas displaying an
347 'IR transformation fingerprint' and areas that do not. In the analysis presented here,
348 it is the predominant fingerprint (i.e. $\geq 50\%$ of datapoints) that was used to classify
349 the lesion as a whole, but, clinically, the worst area rather than the predominant
350 signal may be more important when predicting transformation and selecting
351 appropriate clinical treatment. Future development of this method would investigate
352 how altering the 50% threshold affects the capability of the model, as the sensitivity
353 would be expected to increase as the threshold is reduced but at the expense of
354 specificity and vice versa. This is an important area for clarification in a larger study
355 and is related to clinical needs. Increased sensitivity (i.e. better prediction of
356 transforming status) could lead to decreased surveillance intervals, treatment such
357 as excision or enrolment onto chemoprevention trials. Conversely increased
358 specificity (i.e. better prediction of non-transforming status) might lead to changes
359 in clinical practice to allow for safe discharge or increased follow-up intervals. The
360 identification of areas predicted to have transforming capacity in one NT lesion is
361 more difficult to explain, as none of the NT lesions were excised during the follow-
362 up period. However, it might be hypothesized that very small islands of putatively
363 transforming OED could have been completely excised during the biopsy procedure,
364 unintentionally performing a therapeutic excisional biopsy despite the intention for
365 diagnostic incisional biopsy.

366 Pre-processing of FTIR data is acknowledged to improve the performance of
367 subsequent classification models [33], but the choice of both the protocol and the
368 subsequent modelling method is often highly subjective and its efficacy is dependent
369 on the characteristics of the dataset. Instead of applying pre-processing steps in an
370 arbitrary manner, a novel objective optimization method was used in the current
371 study to maximize the efficiency of the OED transformation classifier by optimizing
372 the parameters of each pre-processing step of the IR data, determining their ideal
373 sequence and identifying the best performing classification method(s), or pipeline,
374 using Bayesian optimization. Considering that the trialed pipelines all contain
375 theoretically sensible pre-processing and classifier combinations, the significant
376 variation in Matthew's correlation coefficient score (-0.53 to 0.37) was surprising.
377 It was noted, however, that many of the negative correlations were obtained when
378 data was not normalized. Infrared data are very sensitive to the amount of substance
379 being probed, so mitigating for inevitable variability in sample thickness and
380 preparation is crucial in order to build models derived from multiple specimens.

381 The decision to use a LOPOCV strategy was taken because the small sample
382 size precluded the division of specimens into larger test/train sets. LOPOCV
383 produces the most robust estimation of the model's true performance, since every
384 combination of patients is used in the training and testing of the model. However,
385 the small sample size (n=17) led to a high standard deviation in the mean sensitivity,
386 specificity and ROC variance at the individual datapoint level because biological
387 variation both between patients and within individual lesions is to be expected and
388 will impinge on the analysis under a LOPOCV strategy. Thus, a larger, multi-center
389 study is required to test the model further.

390 Four of the six wavenumbers attributed the most weight during this
391 classification can be assigned to components of the amide I and II bands [34, 35],
392 which are situated at 1700-1600 cm^{-1} and 1600-1500 cm^{-1} respectively. Absorbance
393 at these wavenumbers can be directly attributed to the vibrating modes of repeating
394 peptide bonds, but the convoluted nature of the amide bands renders the task of
395 attributing specific moieties to particular wavenumbers difficult. The remaining
396 features with the highest weight, and hence discriminating power, are at 1242 cm^{-1}
397 and 1020 cm^{-1} : regions known to be dominated by contributions from DNA, RNA
398 and glycogen [34]. Thus, a relatively high weighting of the IR data obtained at
399 1242 cm^{-1} might be indicative of an increase in DNA aneuploidy in T lesions
400 compared with NT lesions, a known early event in oral carcinogenesis [36, 37].
401 Similarly, the characteristic absorption peak of glycogen is centered at
402 approximately 1030 cm^{-1} and its importance in the experiments presented in this
403 report may correlate with the observation that abundance of the molecule is depleted
404 in pre-malignant tissue as a result of increased proliferation requiring additional
405 energy [38]. This association with glycogen depletion has been applied in the use of
406 Lugol's Iodine staining in an attempt to identify and clear OED at the margins of
407 oral cancer resections [39]. Despite these reassuring correlations between
408 wavenumbers with high weighting in the discriminatory model and previously
409 recognized biological observations, it should be remembered that the absorbance of
410 IR light at any particular wavenumber by a biological tissue is the sum of the
411 absorbance by a number of different biochemical molecules, and it should not be
412 expected that a multivariate analysis combining DNA aneuploidy and glycogen
413 levels will be as effective a discriminator as the model presented here. Future
414 research with larger sample numbers and intention-to-treat biopsy specimens should
415 use multivariate analysis to assess how many key wavenumbers are necessary, in
416 conjunction with clinicopathological variables, to build a clinically useful
417 discriminatory model. This reduction in the number of wavenumbers required for
418 discrimination will, in turn, lead to the development of less expensive IR-based
419 technology, perhaps utilizing quantum cascade lasers (QCLs) [40, 41], that might be
420 employed in routine pathology laboratories.

421

422

423 **Conclusions**

424 This study of a pathologically defined set of OED specimens with known
425 outcome suggests that the analysis of IR data can distinguish lesions with the
426 capacity to transform to oral cancer from those that do not, regardless of OED grade.
427 This represents a novel analysis of FTIR data collected from oral premalignant
428 lesions, as data is commonly correlated with OED stage rather than outcome. The
429 results are encouraging, bearing in mind the small sample size and the inherent
430 clinical and biological limitations of using a small biopsy to reflect a much greater
431 field of potential malignant change, and may come to represent a step forward in the
432 clinical assessment of such lesions that allows improved treatment planning. Further
433 research should concentrate on increasing sample size and complexity to reflect the
434 clinical conundrum and the development of technology to apply the methodology in
435 a timely and cost-effective manner.
436

437 **Acknowledgements**

438 This study was funded by Cancer Research UK C7738/A26196. BGE and CAW
439 were supported by Engineering and Physical Sciences Research Council (EPSRC)
440 PhD studentships.
441

442 **Contributions**

443 BGE: Data curation; Formal analysis; Investigation; Methodology; Resources; Software;
444 Visualization; Writing – original draft; Writing – review & editing.
445 CAW: Formal analysis; Software; Writing – review & editing.
446 AT: Investigation; Resources; Writing – review & editing.
447 PJG: Investigation; Resources; Writing – review & editing.
448 CIS: Investigation; Project administration; Writing – review & editing.
449 SDB: Formal analysis; Funding acquisition; Supervision; Writing – review & editing.
450 PG: Investigation; Resources; Writing – review & editing.
451 RJS: Conceptualization; Funding acquisition; Supervision; Writing – review & editing.
452 PW: Conceptualization; Funding acquisition; Supervision, Writing – review & editing.
453 JMR: Conceptualization; Funding acquisition; Methodology; Project administration;
454 Resources; Supervision; Writing – original draft; Writing – review & editing.

455 **References**

- 456 1. Ferlay J, Ervik M, Lam F, Colombet M, Mery L, Piñeros M, et al. Global
457 Cancer Observatory: Cancer Today. Lyon, France: International Agency for
458 Research on Cancer. 2020 [2nd July 2021]. Available from:
459 [https://gco.iarc.fr/today/data/factsheets/cancers/1-Lip-oral-cavity-fact-](https://gco.iarc.fr/today/data/factsheets/cancers/1-Lip-oral-cavity-fact-sheet.pdf)
460 [sheet.pdf](https://gco.iarc.fr/today/data/factsheets/cancers/1-Lip-oral-cavity-fact-sheet.pdf)
- 461 2. Johnson NW, Jayasekara P, Amarasinghe AA. Squamous cell carcinoma and
462 precursor lesions of the oral cavity: epidemiology and aetiology. *Periodontol*
463 2000. 2011;57(1):19-37. <https://doi.org/10.1111/j.1600-0757.2011.00401.x>
464 PMID: 21781177
- 465 3. Warnakulasuriya S, Kujan O, Aguirre-Urizar JM, Bagan JV, Gonzalez-
466 Moles MA, Kerr AR, et al. Oral potentially malignant disorders: A consensus
467 report from an international seminar on nomenclature and classification,
468 convened by the WHO Collaborating Centre for Oral Cancer. *Oral Dis.*
469 2021;27(8):1862-80. <https://doi.org/10.1111/odi.13704> PMID: 33128420
- 470 4. Barnes L, Eveson JW, Reichart P, Sidransky D. World Health Organization
471 Classification of Tumors. Pathology and Genetics of Head and Neck
472 Tumours. Lyon: IARC Press; 2005.
- 473 5. Panwar A, Lindau R, Wieland A. Management for premalignant lesions of
474 the oral cavity. *Expert Rev Anticancer Ther.* 2014;14(3):349-57.
475 <https://doi.org/10.1586/14737140.2013.842898> PMID: 24559323
- 476 6. Reddi SP, Shafer AT. Oral premalignant lesions: management
477 considerations. *Oral Maxillofac Surg Clin North Am.* 2006;18(4):425-33.
478 <https://doi.org/10.1016/j.coms.2006.08.002> PMID: 18088843
- 479 7. Mehanna HM, Rattay T, Smith J, McConkey CC. Treatment and follow-up
480 of oral dysplasia - a systematic review and meta-analysis. *Head Neck.*
481 2009;31(12):1600-9. <https://doi.org/10.1002/hed.21131> PMID: 19455705
- 482 8. Muller S. Update from the 4th Edition of the World Health Organization of
483 Head and Neck Tumours: Tumours of the Oral Cavity and Mobile Tongue.
484 *Head Neck Pathol.* 2017;11(1):33-40. [https://doi.org/10.1007/s12105-017-](https://doi.org/10.1007/s12105-017-0792-3)
485 [0792-3](https://doi.org/10.1007/s12105-017-0792-3) PMID: 28247230
- 486 9. Kujan O, Khattab A, Oliver RJ, Roberts SA, Thakker N, Sloan P. Why oral
487 histopathology suffers inter-observer variability on grading oral epithelial
488 dysplasia: an attempt to understand the sources of variation. *Oral Oncol.*
489 2007;43(3):224-31. <https://doi.org/10.1016/j.oraloncology.2006.03.009>
490 PMID: 16931119

- 491 10. Warnakulasuriya S, Kovacevic T, Madden P, Coupland VH, Sperandio M,
492 Odell E, et al. Factors predicting malignant transformation in oral potentially
493 malignant disorders among patients accrued over a 10-year period in South
494 East England. *J Oral Pathol Med*. 2011;40(9):677-83.
495 <https://doi.org/10.1111/j.1600-0714.2011.01054.x> PMID: 21762430
- 496 11. Pitiyage G, Tilakaratne WM, Tavassoli M, Warnakulasuriya S. Molecular
497 markers in oral epithelial dysplasia: review. *J Oral Pathol Med*.
498 2009;38(10):737-52. <https://doi.org/10.1111/j.1600-0714.2009.00804.x>
499 PMID: 19903246
- 500 12. Schepman KP, van der Waal I. A proposal for a classification and staging
501 system for oral leukoplakia: a preliminary study. *Eur J Cancer B Oral Oncol*.
502 1995;31B(6):396-8. [https://doi.org/10.1016/0964-1955\(95\)00032-1](https://doi.org/10.1016/0964-1955(95)00032-1) PMID:
503 8746271
- 504 13. Lumerman H, Freedman P, Kerpel S. Oral epithelial dysplasia and the
505 development of invasive squamous cell carcinoma. *Oral Surg Oral Med Oral*
506 *Pathol Oral Radiol Endod*. 1995;79(3):321-9. [https://doi.org/10.1016/s1079-
507 2104\(05\)80226-4](https://doi.org/10.1016/s1079-2104(05)80226-4) PMID: 7621010
- 508 14. Arduino PG, Surace A, Carbone M, Elia A, Massolini G, Gandolfo S, et al.
509 Outcome of oral dysplasia: a retrospective hospital-based study of 207
510 patients with a long follow-up. *J Oral Pathol Med*. 2009;38(6):540-4.
511 <https://doi.org/10.1111/j.1600-0714.2009.00782.x> PMID: 19453839
- 512 15. Ho MW, Risk JM, Woolgar JA, Field EA, Field JK, Steele JC, et al. The
513 clinical determinants of malignant transformation in oral epithelial dysplasia.
514 *Oral Oncol*. 2012;48(10):969-76.
515 <https://doi.org/10.1016/j.oraloncology.2012.04.002> PMID: 22579265
- 516 16. Monteiro L, Mello FW, Warnakulasuriya S. Tissue biomarkers for predicting
517 the risk of oral cancer in patients diagnosed with oral leukoplakia: A
518 systematic review. *Oral Dis*. 2021;27(8):1977-92.
519 <https://doi.org/10.1111/odi.13747> PMID: 33290585
- 520 17. Baker MJ, Byrne HJ, Chalmers J, Gardner P, Goodacre R, Henderson A, et
521 al. Clinical applications of infrared and Raman spectroscopy: state of play
522 and future challenges. *Analyst*. 2018;143(8):1735-57.
523 <https://doi.org/10.1039/C7AN01871A>
- 524 18. Paraskevaidi M, Matthew BJ, Holly BJ, Hugh BJ, Thulya CPV, Loren C, et
525 al. Clinical applications of infrared and Raman spectroscopy in the fields of
526 cancer and infectious diseases. *Applied Spectroscopy Reviews*. 2021;56(8-
527 10):804-68. <https://doi.org/10.1080/05704928.2021.1946076>
- 528 19. Trevisan J, Angelov PP, Carmichael PL, Scott AD, Martin FL. Extracting
529 biological information with computational analysis of Fourier-transform

- 530 infrared (FTIR) biospectroscopy datasets: current practices to future
531 perspectives. *Analyst*. 2012;137(14):3202-15.
532 <https://doi.org/10.1039/C2AN16300D>
- 533 20. Jolliffe IT, Cadima J. Principal component analysis: a review and recent
534 developments. *Philosophical Transactions of the Royal Society A:*
535 *Mathematical, Physical and Engineering Sciences*.
536 2016;374(2065):20150202. <https://doi.org/doi:10.1098/rsta.2015.0202>
- 537 21. Sattlecker M, Stone N, Bessant C. Current trends in machine-learning
538 methods applied to spectroscopic cancer diagnosis. *TrAC Trends in*
539 *Analytical Chemistry*. 2014;59:17-25.
540 <https://doi.org/10.1016/j.trac.2014.02.016>
- 541 22. Baker MJ, Gazi E, Brown MD, Shanks JH, Gardner P, Clarke NW. FTIR-
542 based spectroscopic analysis in the identification of clinically aggressive
543 prostate cancer. *Br J Cancer*. 2008;99(11):1859-66.
544 <https://doi.org/10.1038/sj.bjc.6604753> PMID: 18985044
- 545 23. Pilling M, Gardner P. Fundamental developments in infrared spectroscopic
546 imaging for biomedical applications. *Chemical Society Reviews*.
547 2016;45(7):1935-57. <https://doi.org/10.1039/C5CS00846H>
- 548 24. Ellis BG, Whitley CA, Al Jedani S, Smith CI, Gunning PJ, Harrison P, et al.
549 Insight into metastatic oral cancer tissue from novel analyses using FTIR
550 spectroscopy and aperture IR-SNOM. *Analyst*. 2021;146(15):4895-904.
551 <https://doi.org/10.1039/D1AN00922B>
- 552 25. Li B, Gu ZY, Yan KX, Wen ZN, Zhao ZH, Li LJ, et al. Evaluating oral
553 epithelial dysplasia classification system by near-infrared Raman
554 spectroscopy. *Oncotarget*. 2017;8(44):76257-65.
555 <https://doi.org/10.18632/oncotarget.19343> PMID: 29100309
- 556 26. Banerjee S, Pal M, Chakrabarty J, Petibois C, Paul RR, Giri A, et al. Fourier-
557 transform-infrared-spectroscopy based spectral-biomarker selection towards
558 optimum diagnostic differentiation of oral leukoplakia and cancer. *Anal*
559 *Bioanal Chem*. 2015;407(26):7935-43. <https://doi.org/10.1007/s00216-015-8960-3> PMID: 26342309
- 561 27. Morais CLM, Paraskevaidi M, Cui L, Fullwood NJ, Isabelle M, Lima KMG,
562 et al. Standardization of complex biologically derived spectrochemical
563 datasets. *Nat Protoc*. 2019;14(5):1546-77. <https://doi.org/10.1038/s41596-019-0150-x> PMID: 30953040
- 565 28. Matthews BW. Comparison of the predicted and observed secondary
566 structure of T4 phage lysozyme. *Biochimica et Biophysica Acta (BBA) -*
567 *Protein Structure*. 1975;405(2):442-51. [https://doi.org/10.1016/0005-2795\(75\)90109-9](https://doi.org/10.1016/0005-2795(75)90109-9)

- 569 29. Chicco D, Jurman G. The advantages of the Matthews correlation coefficient
570 (MCC) over F1 score and accuracy in binary classification evaluation. BMC
571 Genomics. 2020;21(1):6. <https://doi.org/10.1186/s12864-019-6413-7> PMID:
572 31898477
- 573 30. Henderson A. ChiToolbox: MATLAB toolbox for handling hyperspectral
574 data generated by SIMS, FTIR and Raman instruments [19/7/2021].
575 Available from: <https://bitbucket.org/AlexHenderson/chitoolbox/>.
- 576 31. Savitzky A, Golay MJE. Smoothing and Differentiation of Data by
577 Simplified Least Squares Procedures. Analytical Chemistry.
578 1964;36(8):1627-39. <https://doi.org/10.1021/ac60214a047>
- 579 32. de Freitas Silva BS, Batista DCR, de Souza Roriz CF, Silva LR, Normando
580 AGC, dos Santos Silva AR, et al. Binary and WHO dysplasia grading
581 systems for the prediction of malignant transformation of oral leukoplakia
582 and erythroplakia: a systematic review and meta-analysis. Clinical Oral
583 Investigations. 2021;25(7):4329-40. [https://doi.org/10.1007/s00784-021-](https://doi.org/10.1007/s00784-021-04008-1)
584 [04008-1](https://doi.org/10.1007/s00784-021-04008-1)
- 585 33. Lasch P. Spectral pre-processing for biomedical vibrational spectroscopy and
586 microspectroscopic imaging. Chemometrics and Intelligent Laboratory
587 Systems. 2012;117:100-14. <https://doi.org/10.1016/j.chemolab.2012.03.011>
- 588 34. Movasaghi Z, Rehman S, ur Rehman DI. Fourier Transform Infrared (FTIR)
589 Spectroscopy of Biological Tissues. Applied Spectroscopy Reviews.
590 2008;43(2):134-79. <https://doi.org/10.1080/05704920701829043>
- 591 35. Barth A, Zscherp C. What vibrations tell us about proteins. Q Rev Biophys.
592 2002;35(4):369-430. <https://doi.org/10.1017/s0033583502003815> PMID:
593 12621861
- 594 36. Donadini A, Maffei M, Cavallero A, Pentenero M, Malacarne D, Di Nallo E,
595 et al. Oral cancer genesis and progression: DNA near-diploid
596 aneuploidization and endoreduplication by high resolution flow cytometry.
597 Cell Oncol. 2010;32(5-6):373-83. <https://doi.org/10.3233/CLO-2010-0525>
598 PMID: 20448331
- 599 37. Alaizari NA, Sperandio M, Odell EW, Peruzzo D, Al-Maweri SA. Meta-
600 analysis of the predictive value of DNA aneuploidy in malignant
601 transformation of oral potentially malignant disorders. J Oral Pathol Med.
602 2018;47(2):97-103. <https://doi.org/10.1111/jop.12603> PMID: 28612463
- 603 38. Aizawa H, Yamada SI, Xiao T, Shimane T, Hayashi K, Qi F, et al.
604 Difference in glycogen metabolism (glycogen synthesis and glycolysis)
605 between normal and dysplastic/malignant oral epithelium. Arch Oral Biol.
606 2017;83:340-7. <https://doi.org/10.1016/j.archoralbio.2017.08.014> PMID:
607 28892665

- 608 39. McCaul JA, Cymerman JA, Hislop S, McConkey C, McMahon J, Mehanna
609 H, et al. LIHNCS - Lugol's iodine in head and neck cancer surgery: a
610 multicentre, randomised controlled trial assessing the effectiveness of Lugol's
611 iodine to assist excision of moderate dysplasia, severe dysplasia and
612 carcinoma in situ at mucosal resection margins of oral and oropharyngeal
613 squamous cell carcinoma: study protocol for a randomised controlled trial.
614 *Trials*. 2013;14:310. <https://doi.org/10.1186/1745-6215-14-310> PMID:
615 24063578
- 616 40. Pilling MJ, Henderson A, Bird B, Brown MD, Clarke NW, Gardner P. High-
617 throughput quantum cascade laser (QCL) spectral histopathology: a practical
618 approach towards clinical translation. *Faraday Discussions*. 2016;187(0):135-
619 54. <https://doi.org/10.1039/C5FD00176E>
- 620 41. Pilling MJ, Henderson A, Gardner P. Quantum Cascade Laser Spectral
621 Histopathology: Breast Cancer Diagnostics Using High Throughput
622 Chemical Imaging. *Analytical Chemistry*. 2017;89(14):7348-55.
623 <https://doi.org/10.1021/acs.analchem.7b00426>
624