# Standardization of complex biologically-derived spectrochemical datasets

Camilo L.M. Morais[a,*,1], Maria Paraskevaidi[a,*,1], Li Cui[d], Nigel J Fullwood[b], Martin Isabelle[c],

Kássio M.G. Lima[e], Pierre L. Martin-Hirsch[f], Hari Sreedhar[h], Júlio Trevisan[g], Michael J

Walsh[h], Dayi Zhang[i], Yong-Guan Zhu[d], Francis L. Martin[a,1]


*[a]School of Pharmacy and Biomedical Sciences, University of Central Lancashire, Preston

PR1 2HE, UK; [b]Division of Biomedical and Life Sciences, Faculty of Health and Medicine,

University of Lancaster, Lancaster LA1 4YQ, UK; [c]Spectroscopy Products Division

Renishaw plc, New Mills, Wotton-under-Edge, Gloucestershire GL12 8JR, UK; [d]Key Lab of

Urban Environment and Health, Institute of Urban Environment, Chinese Academy of

Sciences, Xiamen 361021, China; [e]Institute of Chemistry, Biological Chemistry and

Chemometrics, Federal University of Rio Grande do Norte, Natal 59072-970, Brazil;

[f]Department of Obstetrics and Gynaecology, Lancashire Teaching Hospitals NHS

Foundation, Preston PR2 9HT, UK; [g]Institute of Astronomy, Geophysics and Atmospheric

Sciences, University of São Paulo, Cidade Universitária, R. do Matão, 1226 - Butantã, São

Paulo - SP, 05508-090, Brazil; [h]Department of Pathology, University of Illinois at Chicago,

Chicago, Illinois, USA; [i]School of Environment, Tsinghua University, Beijing 100084, China*




[1]To whom correspondence should be addressed: Email: cdlmedeiros-de-morai@uclan.ac.uk;

Email: mparaskevaidi@uclan.ac.uk; Email: flmartin@uclan.ac.uk; Tel: +44 (0) 1772 89 6482

*Contributed equally

## Abstract

The use of spectroscopy techniques, such as Fourier-transform infrared (FTIR) spectroscopy, has been a successful method to study the interaction of light with biological materials and facilitate novel cell biology analysis. Disease screening and diagnosis, microbiological studies, forensic and environmental investigations make the use of spectrochemical analysis very attractive due to its low cost, minimal sample preparation, non-destructive nature and substantially accurate results. However, there is now an urgent need for repetition and validation of these methods in large-scale studies and across different research groups, which would bring the method closer to clinical and/or industrial, implementation. In order for this to succeed, it is important to eliminate the chance of random spectral alterations caused by inter-individual, inter-instrument and/or inter-laboratory variations. Thus, it is evident that spectral standardization is crucial for the widespread adoption of these spectrochemical technologies. By using calibration transfer procedures, different sources of variations can be normalized into a single model using computational-based methods; therefore, measurements performed under different conditions can eventually generate the same result, eliminating the need for a full recalibration. In this paper, we have constructed a protocol for model standardization using different transfer technologies described for FTIR spectrochemical applications. This is a critical step towards the construction of a practical spectrochemical analysis model for daily routine analysis, where uncertain and random variations are present.

## Introduction

Vibrational spectroscopy has shown great promise as an analytical tool for the investigation of numerous sample types with wide applications in diverse sectors, such as biomedicine, pharmaceutics or environmental sciences. Fourier-transform infrared (FTIR) spectroscopy is one of the preferred techniques for identification of biomolecules through the study of their characteristic vibrational movements. Using chemometric approaches, the system is trained to recognize unique spectral features within a sample, so that when unknown samples are introduced an accurate classification is feasible. Alterations in these measurement parameters could interfere with the spectral signature and produce random variations. Therefore, a crucial step is spectral correction, or standardization, which would provide comparable results and allow system transferability. The idea is that non-biological variations, such as those arising from different users, locations or instruments, will no longer affect the classification result; therefore any collected data could be imported into a central database and handled for further exploration or diagnostic purposes. Several groups and companies worldwide are developing spectrochemical approaches for diagnosis, discrimination and monitoring of diseases, as well as for other uses. Combination of multiple datasets would facilitate the conduction of large-scale studies, which are still lacking in the field of biospectroscopy.

## Sensor-based technologies

Sensor-based technologies are an integral part of daily life ranging from locating sensor-based technology, such as global positioning system (GPS)[1], to image biosensors, such as X-rays[2-5] and γ-rays[6-8], which are used extensively for medical applications. Other powerful approaches that make use of sensor-based technologies toward medical disease examination and diagnostics include circular dichroism (CD) spectroscopy[9-12], ultraviolet

73  (UV) or visible spectroscopy[13,14], fluorescence[15-19], nuclear magnetic resonance (NMR)

74  spectroscopy[20-24] and ultrasound (US) [2,25-28].

75      Over the last two decades, optical biosensors employing vibrational spectroscopy,

76  particularly IR spectroscopy, have seen tremendous progress in biomedical and biological

77  research. A number of studies using the above-mentioned methods have focused on cancer

78  investigation with malignancies such as brain[29-32], breast[33-35], oesophagus[36,37], skin[38-42],

79  colorectal[43-45], lung[46-48], ovarian[49-53], endometrial[50,54,55], cervical[56-59] and prostate[60-63] cancer

80  being some of them. Non-cancerous diseases have also been examined, namely

81  neurodegenerative disorders[64-67], HIV/AIDS[68], diabetes[69-71], rheumatoid arthritis[72,73],

82  cardiovascular diseases[74,75], malaria[76-78], alkaptonuria[79], cystic fibrosis[80], thalassemia[81],

83  prenatal disorders[82,83], macular degeneration[84,85], atherosclerosis[75,86] and osteoarthritis[87-89].

84  Limitations

85      Spectrochemical approaches are advantageous when compared with traditional

86  molecular methods as they provide a holistic status of the sample under interrogation, thus

87  generating typical spectral regions widely known as "fingerprint regions". These methods

88  have also been shown to be rapid, inexpensive and non-destructive while they also improve

89  diagnostic performance and eliminate subjective diagnosis (*e.g*., histopathological diagnosis),

90  where inter- and intra-observer variability are present[90]. However, similarly to any other

91  analytical method, vibrational spectroscopy also comes with some limitations. For instance,

92  prior to FTIR studies, optimization of instrumental settings, sample preparation and operation

93  mode also needs to be conducted in order to improve the spectral quality and molecular

94  sensitivity[91-93]. Overall, the above-mentioned barriers can be overcome after careful

95  consideration of the experimental design.

96      A considerable limitation that is yet under-investigated in the field of spectrochemical

97  techniques is associated with the difficulties entailed in data conformation and system

98    standardization. Currently, there are multiple pilot studies showing promising results but an

99    approach towards standardization for biological applications is lacking. Random variation

100   between studies can originate from differences in instrumentation, operators, and

101   environmental conditions, such as room temperature and humidity.

102          The main objective of this article is to present a protocol for model standardization,

103   which can be applied in FTIR spectrochemical techniques to rule out the chance of random

104   spectral alterations. Inter-individual, inter-instrument, inter-sample and/or inter-laboratory

105   variations can be a source of unwanted, non-biological alterations, thus leading to incorrect

106   conclusions. However, for a method to become reliable and clinically translatable, it is

107   important that measurements performed under different conditions generate comparable

108   results. The aim of the spectral standardization model presented here is to expedite multi-

109   centre studies with large numbers of samples; this would bring these spectrochemical

110   techniques closer to clinical implementation and facilitate life-changing decisions. We

111   describe a protocol that has four main components: (i) sample preparation, (ii) spectral

112   acquisition, (iii) data pre-processing and (iv) model standardization. The current protocol has

113   an in-depth insight obtained from cross-laboratory collaborations with leading experts in the

114   field. This article offers a step-by-step procedure, which can be implemented by a non-

115   specialist in spectrochemical studies. For further information about instrumental and software

116   options, spectral acquisition steps and data analysis for a range of different analytical systems

117   the reader is directed towards additional protocols[91,94-101].

118   Applications

119          Spectrochemical approaches, in combination with computational analysis, have been

120   proven to be effective for biomedical research through facilitating the diagnosis,

121   classification, prognosis, treatment stratification and modulation or monitoring of a disease

122   and treatment. However, these techniques are widely applicable to other fields as well,

123   namely    food    industry[102-105],    toxicology[106-109],    microbiology[110-115],    forensics[116-120],

124   pharmacy[108,121,122], environmental and plant science[123-125], as well as defence and security[126-

125   [128]. Applications of standardization algorithms vary according to the spectral technique and

126   sample matrix studied, where mostly are based on Raman and Fourier-transform near-

127   infrared (FT-NIR) spectroscopy. Table 1 summarizes some standardization applications.

128

129 **Table 1.** Examples of applications involving standardization techniques.

| Sample matrix | Spectroscopic technique | Aim | Ref. |
|---|---|---|---|
| Tissue | Raman | Standardization of various perturbations on Raman spectra for diagnosis of breast cancer based on snap frozen tissues | [129] |
| | Raman | Standardization of spectra acquired in 3 different sites for analysing oesophageal samples based on snap frozen tissues | [130] |
| Cells | Raman | Standardization of spectra acquired with 4 different instruments for classification of three different cultured spore species | [131] |
| Biofluids | FT-NIR | Standardization of spectra acquired with 3 different instruments for measuring haematocrit in the blood of grazing cattle | [132] |
| | LC-MS | Standardization of spectra acquired with 2 different instruments for mapping rendition times and matching metabolite features of subjects diagnosed with small cell lung cancer based on blood serum and plasma samples analysis | [133] |
| Pharmaceutical materials | Raman | Standardization of spectra acquired with 5 different instruments for analysing various pharmaceutical excipients, active pharmaceutical ingredients (APIs) and common contaminants | [134] |
| | FT-NIR | Standardization of spectra acquired with 2 different instruments for simultaneous determination of rifampicin and isoniazid in pharmaceutical formulations | [135] |
| | FT-NIR | Standardization of spectra acquired with 2 different instruments for predicting content of 654 pharmaceutical tablets | [136] |
| Food | FT-NIR | Standardization of spectra acquired with 3 different instruments for predicting parameters in corn samples | [136] [137] |
| | FT-NIR | Standardization of spectra acquired with 2 different instruments for predicting vitamin C in navel orange | [138] |
| | FT-NIR | Standardization of spectra recorded in 4 different labs for determining moisture, proteins and oil content in soy seeds | [139] |
| | FT-NIR | Standardization of spectra acquired by a benchtop and portable instrument for determining total soluble solid contents in single grape berry | [140] |
| | UV-Vis | Standardization of visible spectra acquired with 3 different instruments for measuring pH of Sala mango | [141] |
| Plant | FT-NIR | Standardization of spectra acquired with 2 different instruments for predicting baicalin contents in radix scutellariae samples | [137] |
| | FT-NIR | Standardization of spectra acquired by 2 different instruments and in three physical states (powder, filament and intact leaf) for determining total sugars, reducing sugars and nicotine in tobacco leaf samples | [142] |
| | NMR | Standardization of spectra acquired with 3 different instruments for authenticity control of sunflower lecithin | [143] |
| Cosmetic | CD spectroscopy | Standardization of spectra acquired between standard and real-world samples for determining $Pb^{2+}$ in cosmetic samples | [144] |
| Inorganic substances | FT-IR | Standardization of interferogram spectra acquired with 2 instruments for classifying acetone and $SF_6$ samples | [145] |
| Fuel | FT-IR | Standardization of spectra acquired with 2 different instruments for predicting density of crude oil samples | [146] |

130

131 ## Model transferability

132       Transferability models have been previously developed, however this is still an under-

133 investigated field, especially for biomedical applications. An inclusive standardization

134 protocol that could be implemented in a range of different spectrochemical approaches is of

135 great need. Differences are present even between identical instruments; for instance, changes

136 in signal intensity caused by replacement, alignment or ageing of optical and spectrometer

137 components, natural variations in optics and detectors construction, changes in measurement

138 conditions (temperature and humidity), changes in physical constitution of the sample

139 (particle size and surface texture) and operator discrepancies could all lead to wavenumber

140 shifts and artefacts in the spectra. In all of these cases, prediction errors can become very

141 large, especially when the whole spectrum is used in the model. Standardization techniques

142 aim to generate a uniform spectral response under differing conditions, ensuring the

143 interchangeability of results obtained in different situations, without having to perform a full

144 calibration for each situation.

145       Previous standardization methods include the use of simple slope and bias

146 correction[147,148], direct standardization (DS)[149-153], piecewise direct standardization

147 (PDS)[147,154-156], piecewise linear discriminant analysis (PLDA)[145], guided model

148 reoptimization (GMR)[156], back-propagation neural network (BNN)[145], generalized least

149 squares weighting (GLSW)[157], model updating (MU)[158,159], orthogonal signal correction

150 (OSC)[160,161], orthogonal projections to latent structures (OPLS)[146], wavelet hybrid direct

151 standardization (WHDS)[155], maximum likelihood PCA (MLPCA)[162], Shenk and Westerhaus

152 method (SW)[163,164], positive matrix factorization (PMF)[165,166], artificial neural networks

153 (ANN) drift correction[167], transfer *via* extreme learning machine auto-encoder method

154 (TEAM)[168], calibration transfer based on the maximum margin criterion (CTMMC)[169],

155 calibration transfer based on canonical correlation analysis (CTCCA)[170] and calibration

156 methods, such as wavenumber offset correction, instrument response correction and baseline

157 correction[130].

158 **Direct standardization.** DS is one of the most used methods for data standardization. It was

159 initially proposed to correct relatively large spectral differences between data collected by

160 two instruments[147]. In DS, the entire spectrum from a new secondary response (*e.g.*, a

161 different instrument) is transformed to resemble the spectrum from the primary source (*e.g.*,

162 original instrument)[149]. This is performed based on a linear relationship between the data

163 acquired under different circumstances[158]:

164 $$\mathbf{S}_1 = \mathbf{S}_2 \mathbf{F} \qquad\qquad (01)$$

165 where $\mathbf{S}_1$ represents the data acquired for the primary response; $\mathbf{S}_2$ represents the data

166 acquired for the secondary response; and $\mathbf{F}$ is the transformation matrix that maintains the

167 relationship between $\mathbf{S}_1$ and $\mathbf{S}_2$.

168 The transformation matrix $\mathbf{F}$ is estimated in a least-squares sense by[171]:

169 $$\mathbf{F} = \mathbf{S}_2^+ \mathbf{S}_1 \qquad\qquad (02)$$

170 where $\mathbf{S}_2^+$ is the pseudo-inverse of $\mathbf{S}_2$, calculated by:

171 $$\mathbf{S}_2^+ = (\mathbf{S}_2^T \mathbf{S}_2)^{-1} \mathbf{S}_2^T \qquad\qquad (03)$$

172 in which T stands for the matrix transpose operation.

173 Then, when samples are measured under the secondary system, the signals generated

174 $\mathbf{X}$ are transformed to resemble the primary system response by[158]:

175 $$\hat{\mathbf{X}}^T = \mathbf{X}^T \mathbf{F} \qquad\qquad (04)$$

176 where $\hat{\mathbf{X}}$ is the standardized response for $\mathbf{X}$.

177    Problems related to different background information between instruments can affect

178    the standardization procedure. To correct for this, the standardization process is usually

179    adapted with the background correction method[171], in which the transformation matrix

180    described in Eq. 02 is calculated with a background correction factor ($\mathbf{F}_b$) and an additive

181    background correction vector $\mathbf{b}_s$ as follows:

182    $$\mathbf{S}_1 = \mathbf{S}_2\mathbf{F}_b + \mathbf{1}\mathbf{b}_s^T \tag{05}$$

183    where $\mathbf{1}$ is an all-ones vector and $\mathbf{b}_s$ is obtained by:

184    $$\mathbf{b}_s = \mathbf{s}_{1m} - \mathbf{F}_b^T\mathbf{s}_{2m} \tag{06}$$

185    in which $\mathbf{s}_{1m}$ is the mean vector of $\mathbf{S}_1$ and $\mathbf{s}_{2m}$ is the mean vector of $\mathbf{S}_2$.

186    One of the key steps for DS is the selection of the number of samples to transfer

187    (called "transfer samples"). These are samples from the primary system ($\mathbf{S}_1$) that will be used

188    to transform the signal obtained using the secondary system ($\mathbf{S}_2$). Usually, the procedure for

189    selecting transfer samples is based on sample selection techniques, such as Kennard-Stone

190    (KS) algorithm[172] or leverage[147]. Subsequently, the number of transfer samples is evaluated

191    using a validation set through an arbitrary cost function. For quantification applications, a

192    common cost function is the root-mean-square error of prediction, while for classification one

193    can use the misclassification rate.

194    A disadvantage of DS is that each transformed variable is calculated using the whole

195    spectrum, which carries a high risk of overfitting. The estimation of $\mathbf{F}$ in Eq. (02) is a ill-

196    conditioned problem, because the number of variables may be much larger than the number

197    of standard samples.

198    **Piecewise direct standardization.** PDS is another standardization procedure commonly

199    employed for system transferability. It is based on DS, however it uses windows (*e.g.*,

200   wavenumber portions) to make the standardization process more suitable for smaller regions

201   of the data. When compared to DS, PDS is calculated by using the transformation matrix **F**

202   with most of its off-diagonal elements set to zero[147]. With this, PDS fits minor spectral

203   modifications not covered by DS. PDS is the technique of preference for correcting smaller

204   spectral variations, such as small wavelengths shift, intensity variations, and bands

205   enlargement and reduction[147]. In addition, an advantage of PDS compared to DS is that the

206   local rank of each window will be smaller than the rank of the whole data matrix, which

207   means that the number of standard samples can be smaller, and indeed good results have been

208   obtained with very few samples.

209       One disadvantage of PDS is the need of an additional optimization process, because in

210   addition to the number of transfer samples, PDS also needs a window size optimization,

211   which might lead to a risk of overfitting. Herein, the window size optimization is made using

212   a cost function expressed as the misclassification rate calculated for each window size tested,

213   being evaluated using a validation set where the window with smaller misclassification is

214   selected for final model construction.

215   Experimental Design

216       A specified number of steps are required for a study using vibrational spectroscopy,

217   starting from careful experimental design, protocol optimisation and development of

218   experimental procedure document, sample collection and preparation, spectral collection, pre-

219   processing of the derived information and lastly the use of chemometrics for exploratory,

220   classification and standardization purposes. FTIR spectroscopy is described in more detail in

221   this study, however, the standardization protocol described here can be adapted to a range of

222   techniques, including attenuated total reflection (ATR-FTIR), transmission and transflection

223   FTIR, near-IR (NIR), UV-visible, NMR spectroscopy and MS. Nevertheless, intrinsic

224 features of each technique should be taken into consideration before standardization and the

225 protocol may change depending on the application of interest.

226     A number of biological samples can be analyzed with the above-mentioned analytical

227 methods such as tissues, cytological materials or biological fluids. Sample type and

228 preparation may differ depending on the technique that is employed each time. For instance,

229 IR spectroscopy is limited by water interference at the fingerprint region that can mask the

230 signal of the analyte close to the water peak. This could be addressed with an extra step of

231 sample drying, in contrast to Raman spectroscopy, for example, where water does not

232 generate signal in this region.

233     Typical steps for sample preparation, acquisition of spectra and data pre-processing

234 are briefly presented here. However, the main focus of this protocol is placed on the

235 calibration transfer and standardization procedures. Readers are directed to additional

236 literature for more detailed information regarding sample format and preparation, suitability

237 of substrates, instrumentation settings or available software packages (Table 2) and

238 manufacturers[91,94-96,101,173-176].

239 **Table 2.** Software packages for data standardization.

| Software | Website | Description | Availability |
|---|---|---|---|
| PLS_Toolbox | http://www.eigenvector.com/ | MATLAB toolbox for chemometric analysis. Contains standardization routines using DS, PDS, double window PDS, spectral subspace transformation, GLSW, OSC, and alignment of matrices. | Commercial |
| Unscrambler® X | http://www.camo.com/ | Software for multivariate data analysis and design of experiments. Contains standardization routines using interpolation, bias and slope correction, and PDS. | Commercial |
| OPUS | https://www.bruker.com/ | Spectral acquisition software with data processing features. Contains a standardization routine using PDS. | Commercial |
| Pirouette® | https://infometrix.com/ | Chemometrics modelling software. Contains standardization routines using DS and PDS. | Commercial |

240

Experimental design: sampling

**Sample preparation.** Biological samples have been studied extensively with spectrochemical techniques for disease research. Tissue specimens can be analysed fresh, snap-frozen or formalin-fixed, paraffin-embedded (FFPE). Fresh or snap-frozen histology sections are preferable as they are devoid of contaminants whereas FFPE treatment contributes to characteristic peaks, hindering the biological information. FFPE tissues can be deparaffinized either by chemical methods (*e.g.*, incubation in xylene, hexane or Histo-Clear solutions)[91], which can alter tissue structures and be inefficient for the complete wax removal[177], or by applying chemometrics (*e.g.*, digital dewaxing)[178,179], which keeps the tissue intact but might introduce artefacts due to over- or under-estimation of the wax contribution[177].

Fixatives, such as ethanol, methanol or formalin, are often used for the preservation of cytological material, also generating strong peaks and interfering with the spectra; thus, a washing step is crucial before spectroscopic interrogation. Fixation in tissue or cells for preservation purposes generates protein cross-linking which can cause changes in the spectra, especially on the Amide I peak[180]. Alternatively, cells can be studied live after washing from residual medium.

Preparation and pre-treatment of biological fluids depend on the sample type. Some of the biofluids that have been previously used in spectroscopic studies include blood (whole blood, plasma or serum), urine, sputum, saliva, tears, cerebrospinal fluid (CSF), synovial fluid, ascitic fluid or amniotic fluid[181-183]. A centrifugation step should precede in cases where the cells present in these fluids are not the focus of the study; the supernatant could then be kept for further analysis. In blood-based studies, the user should also consider the anticoagulant of preference (*e.g.*, EDTA, citrate or heparin) as it could generate unwanted spectral peaks[184-186]. Careful planning of experiments as well as consistence throughout a

266    study are of great importance for the generation of robust results. Samples should be very

267    stable, since the spectral differences between the data collected under different situations

268    (*e.g.*, different instruments or temperature) should be directly related to the difference

269    between the systems and not a change caused by chemical or physical degradation of the

270    samples. Optimal sample thickness, suitability of substrates and sample formats can differ

271    from one analytical technique to another and thus the user should decide and tailor these

272    according to the study's objective. Another consideration is the number of freeze-thaw cycles

273    and long-term storage as these could compromise the integrity of the samples[184,187].

274    Preferably, FFPE tissue samples should be analysed after thorough dewaxing and freeze-thaw

275    cycles or long-term storage avoided since these could result in many confounding factors for

276    analysis.

277    **Spectral acquisition.** Depending on the study's objective, FTIR spectral information can be

278    collected using either point spectra or imaging. FTIR spectra can be collected in different

279    operational modes, namely ATR-FTIR, transmission or transflection. Instrument parameters

280    such as resolution, aperture size, interferometer mirror velocity and co-additions have to be

281    optimised before acquisition of spectra to achieve high SNR[91,94]. Metal surfaces can also be

282    used to increase the IR signal in a technique known as surface-enhanced IR absorption

283    (SEIRA)[188,189]. As water interference can mask biological information in IR spectra, the user

284    can purge the spectrometer with dry air or nitrogen gas to reduce the instrument internal

285    humidity, or use computational analysis to remove the water signature. In addition, samples

286    should be dried until all water content evaporates; however, drying of a sample is not without

287    consequences, since chemical changes may occur such as loss of volatile compounds. A

288    background sample is collected regularly to account for any changes in the atmospheric or

289    instrument conditions.

290    For analysing homogenous samples (*e.g.*, biofluids), measurements can be performed

291    by acquiring spectra on different regions of the centre of a drop and across its borders. In

292    transmission measurements, the sample can be measured raw or diluted. Usually, 10 spectra

293    are collected per sample. A higher number of spectral replicas can be performed to decrease

294    the standard-deviation (SD) between measurements, since the SD is proportion to $1/\sqrt{n}$,

295    where $n$ is the number of replicas. For heterogeneously distributed samples (*e.g.*, tissues),

296    spectra should be acquired covering the sample surface as much uniformly as possible, to

297    ensure that all sources of information is contemplated in the spectral data. Samples replicas

298    are also recommended at least as triplicates. For precision estimation, at least six replicates at

299    three levels should be performed. The minimum number of samples for analysis can be

300    estimated using a power test at an 80% power[190]. Further details regarding sampling

301    methodologies for analysing biological materials using FTIR spectroscopy can be found in

302    our previous protocols[91,94].

303    Experimental design: data quality evaluation

304    Before processing, the data can be assessed to identify presence of anomalous

305    behaviours or biased patterns. This can be made initially by visual inspection (e.g.,

306    identification of very anomalous spectra) followed by Hotelling $T^2$ *versus* Q residuals charts

307    using only the mean-centred spectra. PCA residuals[191] can explored to identify biased

308    patterns, in which heteroscedastic distributions are signs of biased experimental

309    measurements; while homoscedastic distributions are associated with good sampling. Also,

310    mistakes performed during experimental data acquisition can be evaluated by $R^2$ values.

311    Negative $R^2$ indicates that the sample variance is smaller than the model residuals variance,

312    which should not happen. SNR can be estimated by dividing the power ($P$) of signal by the

313    power of noise, that is $\text{SNR} = P_{signal}/P_{noise} = \left(A_{signal}/A_{noise}\right)^2$, where $A$ is the amplitude;

314    or by the inverse of the coefficient of variation, when only non-negative variables are

315    measured. Collinearity can be evaluated by calculation of condition number, which is

316    naturally high for spectral data (high collinearity).

317    Experimental design: pre-processing

318        Data pre-processing is employed for maximizing the SNR. This process is

319    fundamental for correcting physical interfering, such as light scattering, different sample

320    thickness, different optical paths and instrumental noise. Therefore, the pre-processing step

321    has fundamental importance to highlight the signal of interest and reduce interfering.

322        For standardization applications, the pre-processing step is also important for

323    reducing differences between the different systems that are used. Before any additional pre-

324    processing, the biofingerprint region should be truncated (*e.g.*, 900-1800 cm$^{-1}$) before

325    analysis. This region contains the main absorptions from biochemical compounds and it

326    suffers minor effects of environmental variability, such as air humidity (free $v$O-H = 3650–

327    3600 cm$^{-1}$, hydrogen-bonded $v$O-H = 3400 – 3300 cm$^{-1}$) and air $CO_2$ ($v_sCO_2$ = 2350 cm$^{-1}$)[192].

328    Table **3** summarizes the main pre-processing techniques for correcting noise in biologically-

329    derived datasets.

**Table 3.** Main pre-processing used for biologically-derived datasets.

| Pre-processing | Interfering | Technique | Advantage | Disadvantage | Optimization |
|---|---|---|---|---|---|
| Savitzky-Golay smoothing[193] | Instrumental noise | ATR-FTIR, FTIR, NIR, Raman, NMR, UV-Vis | Corrects spectral noise without changing the shape of data significantly | The polynomial order and window size for polynomial fit affects the result | The polynomial function should have an order similar to the spectral data (*e.g.*, 2$^{nd}$ order polynomial function for IR data) and the window size should be an odd number and not too small (keeping the noise) or too large (changing the spectral shape) |
| Multiplicative scatter correction (MSC)[194] | Light scattering (Mie scattering), different pressure over the sample when using ATR or probe, different lengths of optical path | ATR-FTIR, FTIR, NIR, Raman, NMR, UV-Vis | Corrects light scattering maintaining the same spectral shape and signal scale | Need of a reference spectrum representative of all measurements | The reference spectrum is regularly set as the average spectrum across all training samples |
| Standard normal variate (SNV)[195] | Light scattering (Mie scattering), different pressure over the sample when using ATR or probe, different lengths of optical path | ATR-FTIR, FTIR, NIR, Raman, NMR, UV-Vis, ~~fluorescence EEM~~ | Corrects light scattering maintaining the same spectral shape | Creates negative signals since the data are centralized to zero (y-scale) | -- |
| Spectral differentiation[193] | Light scattering (Mie scattering), different pressure over the sample when using ATR or probe, different lengths of optical path, background absorption interfering | ATR-FTIR, FTIR, NIR, Raman, NMR, UV-Vis | Corrects light scattering and baseline problems; highlights smaller spectral differences | Changes the signal scale, shifts the data and increases noise | The order of the derivative function should be used carefully to avoid increased noise (usually 1$^{st}$ or 2$^{nd}$ order differentiation is preferred). The differentiation can be coupled to Savitzky-Golay smoothing |
| Baseline correction[196] | Background absorption interfering | ATR-FTIR, FTIR, NIR, Raman, NMR, UV-Vis, MS | Corrects the baseline maintaining the same spectral shape | -- | There are many methods for baseline correction (*e.g.*, rubber band, automatic weighted least squares, Whittaker filter). The method chosen should be maintained consistent for all systems used |
| Normalization[90] | Different sample thickness and concentration | ATR-FTIR, FTIR, Raman | Avoids influence of non-desired signals among the | The normalization might hide signal differences | -- |

| samples | between samples at important bands, such as Amide I and Amide II; and also may introduce non-linearities |
| --- | --- |

331    Figure 1 depicts the effect of a pre-processing approach employed for a blood plasma

332    dataset acquired under different experimental conditions (*i.e.*, different systems and

333    operators). In this Figure, the reduction of the spectral differences between the systems is

334    evident after data pre-processing (Savitzky-Golay smoothing, MSC, baseline correction and

335    normalization).



336

337    **Figure 1.** Average (a) raw and (b) pre-processed IR spectra for healthy control samples

338    across three different systems (A, B and C). Average (c) raw and (d) pre-processed IR spectra

339    for healthy control samples across two different operators (Operator 1 and 2).

340

341    After the pre-processing techniques displayed in Table 3, scaling should be employed

342    as most classification methods require all the variables (*e.g.*, wavenumbers) in the dataset to

343    be at the same scale in order to work properly.

344    For spectral data, mean-centring (also referred as "standardization" by Hastie et al.[197])

345    is a very reasonable approach, after which all variables in the dataset will have zero mean.

346    When data contain values represented by different scales (*e.g.*, after data fusion using both IR

347    and Raman spectra), block-scaling should be used, where each block of data would have the

348    same sum-of-squares (normally after mean-centring).

349    Another important aspect of pre-processing is the order in which each step is applied.

350    Pre-processing should be employed in a logical order so that the next pre-processing step is

351    not affected by the previous one. For example, pure spectral differentiation cannot be

352    employed before smoothing, since the spectral differentiation will increase the original noise.

353    Therefore, smoothing should be applied before differentiation. Albeit, Savitzky-Golay routine

354    incorporates smoothing and spectral differentiation so, in practical terms, these can be

355    performed together. To summarise, the suggested order of pre-processing is as follows:

356    1.  Spectral Truncation

357    2.  Smoothing

358    3.  Light scattering correction

359    4.  Baseline correction

360    5.  Normalization

361    6.  Scaling

362    When using different instruments but same type of sample, the pre-processing steps

363    should be the same for the data acquired under different circumstances.

Experimental design: data analysis

365 **Sample splitting.** Sample splitting is fundamental for constructing a predictive chemometric

366 model. The splitting procedure can be performed manually or by computer-based

367 methodologies. Manual splitting can generated biased results, therefore computational-based

368 split is more recommended. In this case, some strategies includes random selection,

369 leverage[147] or the KS algorithm[172]. KS works based on Euclidian distance calculation by

370 firstly assigning the sample with the maximum distance to all other samples to the calibration

371 set, and then by selecting the samples which are as far away as possible from the selected

372 samples to this set, until the designed number of selected samples is reached. This ensures

373 that the calibration model will contain samples that uniformly cover the complete sample

374 space, where no or minimal extrapolation of the remaining samples are necessary; avoiding

375 problems of manual or random selection, such as non-reproducibility and non-representative

376 selection. Usually, the dataset is split with 70% of the samples assigned for training, 15% for

377 validation and 15% for test. In this case, the test set is dependent on the initial group of

378 samples measured, and it is not a regular independent test set where a new set of similar

379 samples are measured.

380 **Exploratory analysis.** Exploratory analysis is an important tool to provide an initial

381 assessment of the data. Using exploratory analysis, the analyst can see the clustering patterns

382 and then draw conclusions related to the nature of samples, outliers and experimental errors.

383 One of the most common techniques for exploratory analysis is principal component analysis

384 (PCA), in which the original data are decomposed into a few principal components (PCs)

385 responsible for most of the variance within the original dataset. The PCs are orthogonal to

386 each other and are generated in a decreasing order of explained variance, so that the first PC

387 represents most of the original data variance, followed by the second PC and so on[198].

388 Mathematically the decomposition takes the form:

389  $$X = TP^T + E \qquad\qquad\qquad\qquad\qquad\qquad (07)$$

390  where **X** represents the pre-processed data (*e.g.*, pre-processed samples' spectra); **T** are the

391  scores; **P** are the loadings; and **E** are the residuals.

392  The PCA scores represent the variance in the sample direction and they are used to

393  assess similarities/dissimilarities among the samples, thus detecting clustering patterns. The

394  PCA loadings represent the variance in the variable (*e.g.*, wavenumber) direction and they are

395  used to detect which variables show the highest importance for the pattern observed on the

396  scores. The PCA loadings are commonly employed as a tool for searching spectral markers

397  that distinguish different biological classes[199]. The PCA residuals represent the difference

398  between the decomposed and original data and can be used to identify experimental errors.

399  Ideally, the PCA residuals should be random and close to zero, representing a heteroscedastic

400  distribution. Otherwise, they can indicate experimental bias according to a homoscedastic

401  distribution.

402  For standardization applications, PCA is a fast, intuitive and reliable tool to observe if

403  there are differences between the spectra acquired by different systems. Ideally, if the same

404  sample is measured under different conditions (different laboratories, instrument

405  manufacturers or user operators) their PCA scores should be random and completely

406  superposed. If a discrimination pattern is observed on the PCA scores, then it is indicative

407  that the data need standardization. Figure 2 illustrates a PCA scores plot from the same

408  samples (blood plasma of healthy controls) measured using three IR instruments before (Fig.

409  2a) and after (Fig. 2b) DS. Even though the samples in Fig. 2a are pre-processed, three

410  different clusters are still evident. After DS the samples measured using different systems are

411  normalized into a single cluster.

412

**Figure 2.** (a) PCA scores for healthy control samples across three different instruments (A, B and C) after pre-processing but before DS; (b) PCA scores for healthy control samples across three different instruments (A, B and C) after DS. The dotted blue circle shows 95 % confidence ellipse (two-sided).

**Outlier detection.** Outlier detection is important to prevent samples, which differ from the original dataset, from affecting the results using predictive models. Outliers can be attributed to experimental errors, such as inconsistent sample preparation or spectral acquisition, or to larger experimental noise, such as Johnson noise, shot noise, flicker noise and environmental noise. These samples can have large leverage for classification, masking the real signal from the samples of interest; therefore, it is advised that they be removed from the dataset used to train the predictive model.

To detect outliers, techniques such as Jack-knife[200], Z-score[201] or *K*-modes clustering[202] can be utilised among others[203]. One of the most popular and visually intuitive technique for detecting outliers is the Hotelling $T^2$ vs Q residual test[204]. In this test, a chart is created using the Hotelling $T^2$ values in x-axis and the Q residuals in the y-axis, generating a scatter plot. The Hotelling $T^2$ represents the sum of the normalized squared scores, which is

23

430 the distance from the multivariate mean to the projection of the sample onto the PCs[205]. The

431 Q residuals represent the sum of squares of each sample in the error matrix, thus measuring

432 the residues between a sample and its projection onto the PCs[205]. All samples far from the

433 origin of this graph are considered outliers and should be removed one at a time, as the PCA

434 is highly influenced by the samples that are included in the model. Samples with high values

435 in both Hotelling $T^2$ and Q residuals are the worst outliers; while samples with high values in

436 only one of these axis are the second worst outliers. Supplementary Method 1 illustrates an

437 example for outlier detection. Squared confidence limits can be draw based on this graph;

438 however, this can hinder outlier detection. For example, in squared confidence limits at a

439 95% level, certain amount of data-points (5%) are set outside these limits.

440 **Classification.** Classification techniques are employed for sample discrimination. Using

441 chemometric analysis, one can distinguish classes of samples based on their spectral features

442 and then make further predictions based on these. The prediction capability of a classification

443 model should be evaluated with external samples (unknown samples) through the calculation

444 of figures of merit, including accuracy (proportion of samples correctly classified considering

445 true positives and true negatives), sensitivity (proportion of positives that are correctly

446 identified) and specificity (proportion of negatives that are correctly identified)[206].

447 There are many types of classification techniques for spectral data. Table 4

448 summarizes the main classification techniques employed for biospectroscopy applications,

449 along with their advantages and disadvantages.

450

451

452

453

454 **Table 4.** Classification techniques.

| Classification Technique | Advantage | Disadvantage |
|---|---|---|
| Linear discriminant analysis (LDA)[207] | Simplicity, fast calculation | Needs data reduction, does not account for classes having different variance structures, greatly affected by classes having different sizes |
| Quadratic discriminant analysis (QDA)[207] | Fast calculation, accounts for classes having different variance structures, not much affected by classes having different sizes | Needs data reduction, higher risk of overfitting |
| Partial least squares discriminant analysis (PLS-DA)[208] | Fast calculation, high accuracy | Greatly affected by classes having different sizes, needs optimization of the number of latent variables (LVs) |
| K-Nearest Neighbours (KNN)[209] | Simplicity, non-parametric, suitable for large datasets | Time consuming, needs optimization of the distance calculation method and $k$ value, highly sensitive to the "curse of dimensionality"[197] |
| Support vector machines (SVM)[210] | Non-linear classification nature, high accuracy | High complexity, high risk of overfitting, needs optimization of kernel function and SVM parameters, time consuming |
| Artificial neural networks (ANN)[211] | Non-linear classification nature, ability to work with incomplete knowledge, high accuracy | High computational cost, needs optimization of the number of neurons and layers, no interpretability ("black box" model) |
| Random forests[212] | Non-linear classification nature, high accuracy, relatively low computational cost | High risk of overfitting, needs optimization of the number of trees, no interpretability ("black box" model) |
| Deep learning approaches[213] | Non-linear classification nature, native feature extraction (e.g., in convolutional neural networks (CNN)), local spatial coherence (CNN), high accuracy | High computational cost, needs hyperparameter optimization, needs large datasets, time consuming, no interpretability ("black box" model) |

455

456     When employing classification techniques, one must follow a parsimony order[214],

457 where the simplest algorithms should be used first, reducing the need for more complex

458 algorithms which would require more optimization steps. An order for using these

459 classification algorithms is: LDA>PLS-DA>QDA>KNN>SVM>ANN>Random

460 forests>Deep learning approaches, from the simplest to the most complex.

461     Classification algorithms can be coupled to feature extraction and feature selection

462 techniques in order to reduce data collinearity/redundancy, thus reducing the risk of

463 overfitting in the classifier training, and speeding up such training, as there are less variables

464 involved. An additional benefit of such a feature extraction/selection step is to provide
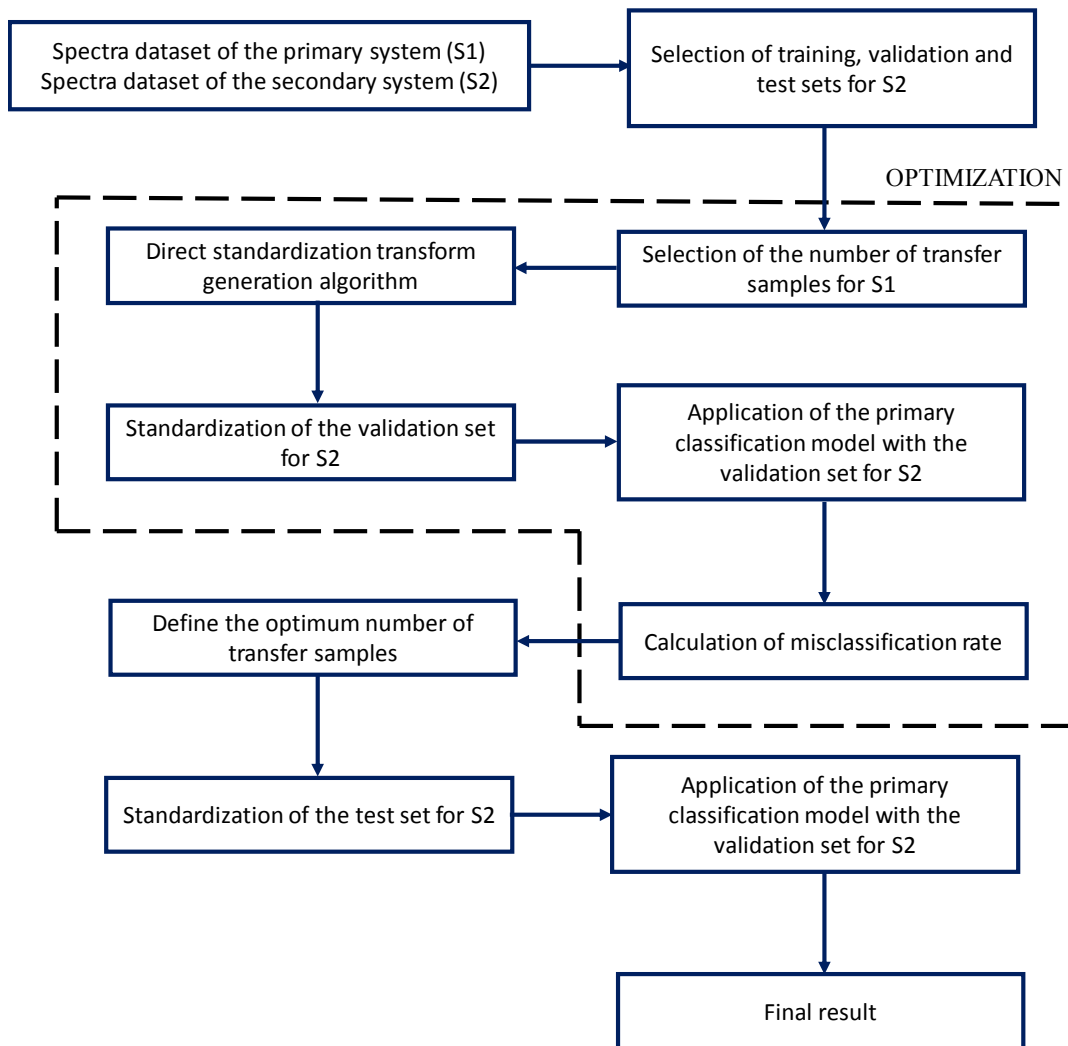
465  spectral    markers    identification    as    a    "side-effect"    (depending    on    the    feature

466  extraction/selection method applied). For feature extraction, the most popular technique is

467  PCA. In this case, a PCA is firstly applied to the data, and then the PCA scores are used as

468  the input variables (instead of the wavenumbers data points) for the classification techniques

469  mentioned above[215]. PLS-DA is also a feature extraction technique[208], and normally it

470  performs better than a PCA followed by LDA, as the scores from a PCA does not necessarily

471  describe the difference between the samples, but rather the variance in the data. In PLS-DA, a

472  partial least squares (PLS) model is applied to the data in an interactive process reducing the

473  original variables to a few number of LVs, where a LDA is used for classifying the groups[216].

474  Other discriminant classifiers, in particular QDA, also could be used in this classification step

475  to circumvent problems observed with LDA. For feature selection, there are many techniques

476  commonly employed in biological datasets, including genetic algorithm (GA)[217] and

477  successive projections algorithm (SPA)[218]. The variables (*e.g.*, wavenumbers) selected by

478  these techniques are used as input variables for the classification models described in Table 2.

479  An important advantage of GA is its relatively low-computational cost compared to SPA and

480  reduction of data collinearity. Furthermore, GA-based techniques are intuitive and simple to

481  understand in the algorithmic sense but they also have a non-deterministic nature and require

482  optimization of many parameters. SPA's advantage relies on its deterministic nature, minor

483  parameter optimization and reduction of data collinearity, however, it is very time

484  consuming. For hyperspectral imaging, feature selection also can be performed by Minimum

485  Redundancy Maximum Relevance (mRMR) algorithm[219], where the selection process is

486  based on maximizing the relevance of extracted features and simultaneously minimize

487  redundancy between them.

488  **Standardization.** Data standardization should be employed when a primary classification

489  model is built and new data comes to be predicted from a secondary system (different

490  laboratory or instrument manufacturers), or when there is a change in instrument components

491  (*e.g*., laser, gratings, etc.) or when the data of the chemometric model are acquired under

492  different circumstances (different analysts, days, instrumental settings, etc.). As previously

493  mentioned, the most common and reliable methods for data standardization are the DS and

494  PDS algorithms. These methods can be found in a few software packages (described in Table

495  3).

496  Figure 3 summarises the standardization protocol using DS applied to spectra

497  acquired under different conditions. The first step consists of applying KS algorithm for

498  selecting the number of transfer samples from the primary system as well as the number of

499  training samples for the secondary systems, which is ideally 70% of the dataset. Thereafter,

500  the DS transform generation algorithm is employed to estimate the transform matrix. The

501  validation set of the secondary system is then used with the classification model of the

502  primary system to evaluate the optimum number of transfer samples. This optimization step

503  is repeated depending on the number of transfer samples from the primary system. After this

504  number is defined, the validation set of the secondary system is finally standardized and the

505  final classification model is subsequently applied. This procedure is realized with a certain

506  number of samples measured in all instruments being standardized. This procedure should be

507  realized in as similar manner as possible to reduce spectral differences. After the model is

508  standardized and proper validated, new external samples can be measured in any of the

509  instruments and predicted by the standardized classification model.

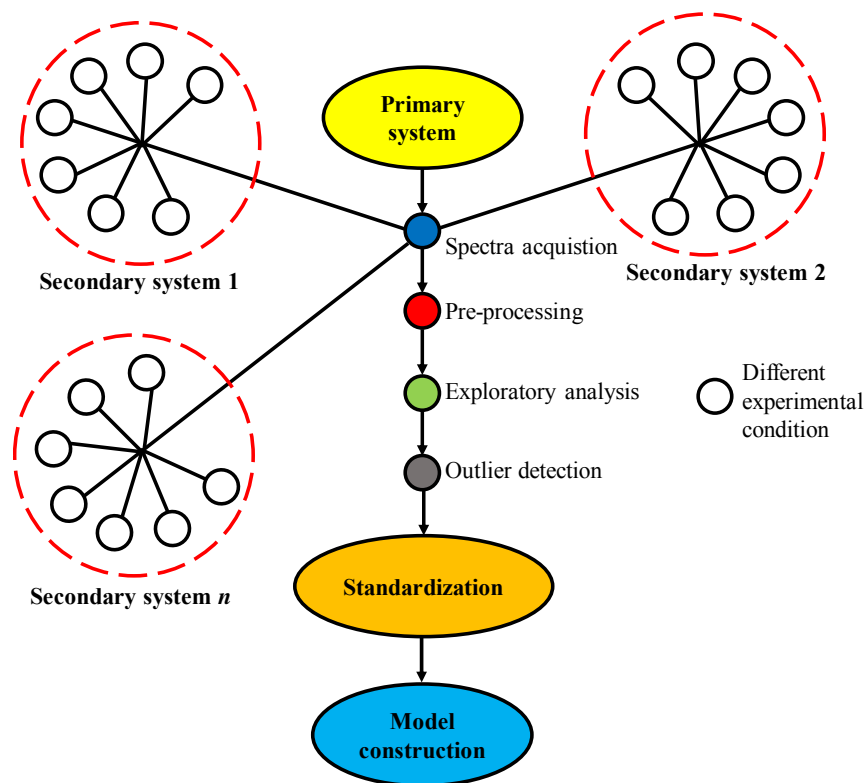**Figure 3.** Flowchart for standardization using Direct Standardization (DS).

For PDS, an extra step is added after defining the number of transfer samples to estimate the optimum window size. The dashed region in Fig. 3 is repeated according to the window size.

For multi-laboratory studies the flowchart depicted in Fig. 4 illustrates how the standardization protocol should be employed.

518



519 **Figure 4.** Flowchart for a standardization protocol using different experimental conditions.

520       In Fig. 4, spectra acquired under different experimental conditions are used for a

521 global standardization model. A primary system should be designated and then all spectra

522 from secondary systems are equally pre-processed, followed by an exploratory analysis to

523 assess samples' similarities/dissimilarities, outlier detection, standardization by the method

524 depicted in Figure 3; the final model construction follows last. With this, all sources of

525 variations present in different systems can be included into a general chemometric model.

526

527 MATERIALS

528 REAGENTS

529 • Biological samples (tissue, cells, biofluids).

530    ▲ **CRITICAL** Human samples should be collected with appropriate local institutional

531    review board for ethical approval and adhere to the Declaration of Helsinki principles.

532    Similarly, for studies involving animals, all experiments should be performed in

533    accordance with relevant guidelines and regulations. Ethical approval has to be obtained

534    before any sample collection.

535    • Optimal cutting temperature (OCT) compound (Agar Scientific, cat. no. AGR1180)

536    • Liquid nitrogen (BOC, CAS no. 7727-37-9) **! CAUTION** Asphyxiation hazard; make

537    sure room is well ventilated. Causes burns; wear face shield, gloves and protective

538    clothing.

539    • Paraplast Plus paraffin wax (Thermo Fisher Scientific, cat. no. SKU502004)

540    • Isopentane (Fisher Scientific, cat. no. P/1030/08) **! CAUTION** Extremely flammable,

541    irritant, aspiration hazard and toxic; use in a fume hood.

542    • Distilled water

543    • PBS (10×; MP Biomedicals, cat. no. 0919610)

544    • Virkon (Antec, DuPont, cat. no. A00960632)

545    • Trypsin–EDTA (0.05%, Sigma-Aldrich, Thermo Fisher Scientific cat. no. 25300054)

546

547    **Anticoagulants**

548    • EDTA (Thermo Fisher Scientific, BD Vacutainer, cat. no. 02-687-107 )

549    • Sodium citrate (Thermo Fisher Scientific, BD Vacutainer)

550    • Lithium/sodium heparin (Thermo Fisher Scientific, BD Vacutainer)

551

552    **Fixative and preservative agents**

553    • Formalin, 10% (vol/vol; Sigma-Aldrich, cat. no. HT501128) **! CAUTION** Potential

554    carcinogen, irritant and allergenic; use in a fume hood.

555 • Ethanol (Fisher Scientific, cat. no. E/0600DF/17)

556 • Methanol (Fisher Scientific, cat. no. A456-212) **! CAUTION** Toxic vapours; use in a

557     fume hood.

558 • Acetone (Fisher Scientific, cat. no. A19-1) **! CAUTION** Acetone vapors may cause

559     dizziness; use in a fume hood.

560 • ThinPrep (PreservCyt Solution, Cytyc Corp)

561 • SurePath (Becton Dickinson Diagnostics)

562

563 **Dewaxing agents**

564 • Xylene (Sigma-Aldrich, cat. no. 534056) **! CAUTION** Potential carcinogen, irritant and

565     allergenic; use in a fume hood.

566 • Histo-Clear (Fisher Scientific, cat. no. HIS-010-010S) **! CAUTION** It is an irritant.

567 • Hexane (Fisher Scientific, cat. no. 10764371) **! CAUTION** Extremely flammable liquid,

568     can cause skin irritation; use protective equipment as required; use in a fume hood.

569

570 EQUIPMENT

571 • Microtome (Thermo Fisher Scientific, cat. no. 902100A; or cat. no. 956651)

572 • Wax dispenser (Electrothermal, cat. no. MH8523B)

573 • Sectioning bath (Electrothermal, cat. no. MH8517)

574 • Centrifuge (Thermo Fisher Scientific, cat. no. 75002410)

575 • Desiccator (Thermo Fisher Scientific, cat. no. 5311-0250)

576 • Desiccant (Sigma-Aldrich, cat. no. 13767)

577 • Laser power meter (Coherent, cat. no. 1098293)

578 • Spectrometer

579 • Computer system

580    **Substrates**

581    ▲ **CRITICAL** Substrate should be carefully chosen depending on the spectrochemical

582    approach that will be used.

583    • Low-E slides (Kevley Technologies, CFR)

584    • $BaF_2$ slides (Photox Optical Systems)

585    • $CaF_2$ slides (Crystran, cat. no. CAFP10-10-1)

586    • Silicon multi-well plate (Bruker Optics)

587    • Glass slides (Fisher Scientific, cat. no. 12657956)

588    • Quartz slides (UQG Optics, cat. no. FQM-2521)

589    • Aluminum-coated slides (EMF, cat. no. AL134)

590    • Mirrored stainless steel (Renishaw, cat. no. A-9859-1825-01)

591

592    REAGENT SETUP

593    **Tissue** For FFPE tissue, the excised specimen is immersed in fixative (*e.g.*, formalin),

594    dehydrated in ethanol, cleared in xylene and embedded in paraffin wax. Specimens can then

595    be stored indefinitely at room temperature. For snap-frozen tissue, the specimen is immersed

596    in OCT, followed by cooling of isopentane with liquid $N_2$.

597    ▲ **CRITICAL** Snap-frozen tissue should be thawed before analysis. Spectroscopic analysis

598    should be performed directly after excision in case of fresh tissue to avoid sample

599    degradation.

600    **Cells** Cells can be treated with a suitable fixative or preservative solution or studied alive.

601    ▲ **CRITICAL** In case cells are fixed or stored in a preservative solution, a number of

602    washing steps using centrifugation should be followed prior to spectroscopic analysis to

603    remove unwanted signature. If cells are studied alive, optimum living conditions (*e.g.*, growth

604     medium, temperature and pH) should be maintained; washing of live cells from medium is

605     also necessary.

606     **Biofluids** Biofluids can be collected in designated, sterile tubes using standard operating

607     procedures to achieve uniformity of performance. Preparation of biofluids depends on the

608     sample type and the experiment's objective. If cellular material is not directly studied, it

609     should be removed from the biofluid before storage. Biofluids can be analysed right after

610     their collection or stored at a -80°C freezer.

611     ▲ **CRITICAL** If biofluids have been stored in a freezer, it is essential that they are fully

612     thawed before acquiring aliquots for spectroscopic analysis.

613     ▲ **CRITICAL** Users are advised to store biofluids in smaller, single-use aliquots at -80°C to

614     avoid repeated freeze-thaw cycles.

615

616     EQUIPMENT SETUP

617     The user can choose from a range of different instrumental setups and spectral acquisition

618     modes. General information about FTIR systems is provided below. For more details about

619     equipment setup see refs.[91,94,95].

620     The FTIR spectrometer can be left on for long periods of time. Before spectral acquisition,

621     the user should check the interferogram signal for amplitude and position and keep a record

622     of the measurements.

623     ▲ **CRITICAL** For detectors that require a prior cooling step using liquid nitrogen (*e.g.*,

624     mercury cadmium telluride (MCT) detectors), the signal should be allowed to stabilize for

625     approximately 10 min before data collection.

626 ▲ **CRITICAL** In case that the interferogram signal deviates from the last measurement, re-

627 alignment or part replacement may be required.

628 **Software**: Software for spectral acquisition is typically provided by the manufacturer.

629 Software packages for spectral analysis and data standardization are provided in Table 3.

630 ## PROCEDURE

631 ### Sample preparation

632 **1|** Prepare the biological samples for spectrochemical analysis using the following steps:

633 option A for FFPE tissue samples, option B for snap-frozen or fresh tissue samples, option C

634 for cells and option D for biofluids.

635 ▲ **CRITICAL** Sample preparation is briefly presented in this protocol. More details about

636 sample preparation can be found in ref.[91,94,95].

637 **(A) Tissue (FFPE)** ● **TIMING** **1-1.5 h**

638 (i) Acquisition of FFPE tissue blocks.

639 (ii) Whole tissue block has to be sectioned using a microtome to obtain tissue sections

640 at desired thickness (2-10 μm).

641 ▲ **CRITICAL** Cooling of the tissue on an ice block allows easier sectioning.

642 (iii) Tissue ribbons are floated in a warm $H_2O$ bath and then deposited onto the

643 substrate of choice.

644 (iv) The tissue slide is then allowed to dry either at room temperature (30 min) or in a

645 60°C oven (10 min).

646 ▲ **CRITICAL** The tissue slide may be dried in the oven for longer periods of time,

647 depending on the type of tissue, to ensure optimal melting of the wax initially.

648       (v) Dewaxing is then performed by three sequential immersions in a dewaxing reagent

649    such as fresh xylene, Histo-Clear solution or hexane (at least 5 min).

650    ▲ **CRITICAL** Thorough dewaxing is important for eliminating all spectral peaks attributed

651    to paraffin.

652       (vi) Tissue slide is immersed in acetone or ethanol (5 min) to remove the xylene and

653    then left to air-dry.

654    ■ **PAUSE POINT** Slides can be stored in a desiccator at room temperature for at least 1

655    year.

656    **(B) Tissue (Snap-frozen or fresh)** ● **TIMING** **2 h + drying time (3 h for FTIR only)**

657    ▲ **CRITICAL** Snap-frozen tissue can be stored at -80°C for several months.

658    ▲ **CRITICAL** For fresh tissue, proceed to step 1B(iii).

659       (i) Acquire snap-frozen tissue from freezer and place onto a cryostat (30 min) to allow

660    the tissue to reach the cryostat's temperature (-20°C).

661       (ii) Tissue block can be sectioned using the cryostat to obtain tissue sections at desired

662    thickness (8-10 μm).

663       (iii) The tissue sections are deposited onto an appropriate substrate before spectra are

664    collected.

665    ▲ **CRITICAL** For FTIR studies the tissue sections need to dry for at least 3 h to remove the

666    $H_2O$ interference with the IR spectra.

667    ▲ **CRITICAL** Exposure to light should be minimised to prevent sample degradation due to

668    oxidation.

**(C) Cells (fixed or live)** ● **TIMING 30 min + desiccation time (3 h for FTIR only)**

▲ **CRITICAL** If cells are studied live proceed to step 1C(ii)

     (i) Fixed cells need to be washed from the fixative or preservative solution to remove any spectral interference in the fingerprint region. Three sequential washes with distilled $H_2O$ or PBS have been shown to remove unwanted peaks.

     (ii) Live cells in suspension have to be detached from the growth substrate using trypsin and then washed from the medium and trypsin with PBS (×3 times).

▲ **CRITICAL** All reagents should be warmed to 37°C to reduce the shock to cells and maintain morphology.

     (iii) After the final wash, the remaining cell pellet is resuspended in distilled $H_2O$ and mounted on a substrate of choice.

▲ **CRITICAL** The final suspension of cells should be evenly deposited on the slide either by cytospinning or by micro-pipetting.

▲ **CRITICAL** For FTIR studies the sample needs to dry for at least 3 h.

**(D) Biofluids (frozen or fresh)** ● **TIMING 5 min + thawing (20 min) + drying (1-1.5 h)**

▲ **CRITICAL** If biofluids are analysed fresh, immediately after collection, continue to step 1D(ii).

     (i) Acquire biofluids from the -80°C freezer and allow them to fully thaw.

     (ii) Mix or gently vortex the sample before obtaining the desired volume for analysis.

▲ **CRITICAL** Only a small amount of the biofluid is typically required for spectroscopic studies (1-100 μL). However, this depends and should be tailored according to the study and experimental design.

691        (iii) Deposit the biological fluid onto an appropriate substrate.

692        ▲ **CRITICAL** For ATR-FTIR spectroscopic studies, an alternative option is to deposit the

693        sample directly on the ATR crystal instead of a substrate if the instrumentation setting allows

694        (*i.e.*, if crystal is facing upwards). However, if the sample is sufficiently thick (>2-3 μm) to

695        avoid substrate interference, then the use of a holding substrate is advantageous as it allows

696        measurements from multiple locations as well as longer storage.

697        ▲ **CRITICAL** For FTIR studies the sample needs to dry adequately before spectroscopic

698        analysis (50 μl dry within approximately 1 h at room temperature). Drying can be sped up by

699        using a gentle stream of air.

700        Spectral acquisition

701        **2|** Spectrochemical information can be collected as follows for FTIR spectroscopy.

702        ▲ **CRITICAL** Spectral acquisition is briefly presented in this protocol. More details can be

703        found in ref.[91,94,95].

704        **FTIR spectroscopy** ● **TIMING 2 - 5 min per spectrum**

705        (i) Settings should be optimised before a new study to increase the SNR (see

706        'Experimental: spectral acquisition').

707        ▲ **CRITICAL** Some of the parameters that need to be adjusted include the resolution,

708        spectral region of interest, co-additions, aperture size, interferometer mirror velocity, and

709        interferogram zero-filling.

710        (ii) Depending on the sampling mode that has been chosen (ATR-FTIR, transmission

711        or transflection), sample is deposited onto the appropriate holding substrate.

712        (iii) Load the sample and visualise the region of interest; information can then be

713        acquired either as point map or as image maps.

714    ▲ **CRITICAL** Typically, 5-25 point spectra are collected per sample while for image maps

715    the step size should be the same or smaller than the selected aperture size divided by two.

716    Sampling can be performed with 6 replicates in 3 levels.

717    ▲ **CRITICAL** A background spectrum should be acquired before every sample to account

718    for atmospheric changes.

719    ▲ **CRITICAL** To improve reproducibility and decrease differences between the data

720    collected by different operators, the spectral resolution should be set constant, since it can

721    cause major differences between data collected across different experimental setups.

722    ▲ **CRITICAL** The pressure applied on the sample in the ATR mode affects the signal

723    intensity (*i.e.*, absorbance) between data collected by different instruments and operators.

724    Thus, the pressure applied on the sample should be as closest as possible across different

725    experimental setups to reduce differences between the spectra collected.

726    ■ **PAUSE POINT** Save the acquired data in a database until further analysis.

727    Data quality evaluation ● **TIMING 15 min – 4 h (depending on the size of the dataset)**

728    ▲ **CRITICAL** Before pre-processing, the raw data can be evaluated using some quality tests

729    to identify anomalous spectra or biased patterns. This can be made by visual inspection of the

730    collected spectra followed by Hotelling $T^2$ *versus* Q residuals charts using only the mean-

731    centred data, and analysis of PCA residuals.

732    Data pre-processing ● **TIMING 15 min – 4 h (depending on the size of the dataset)**

733    ▲ **CRITICAL** Steps 1-6 below can vary depending on the nature of the dataset. Table 1

734    provides more details about these pre-processing steps. In case of an ATR-FTIR dataset

735    acquired under different experimental conditions, the pre-processing method should follow

736    this order:

737    **1. Cutting at biofingerprint region (900-1800 cm$^{-1}$).** The spectra should be truncated

738        to the biofingerprint region to reduce atmospheric interference.

739    **2. Savitzky-Golay smoothing for removing spectral-noise.** Window size varies

740        according to the size of the spectra dataset (*e.g.*, wavenumber). The window size

741        should be an odd number and the analyst should vary it from 3 to 21 and observe how

742        the spectra change (in shape) and how the noise is reduced. The smallest window that

743        removes the noise considerably whilst maintaining the original spectral shape should

744        be used. Using a spectral resolution of 4 cm$^{-1}$, the biofingerprint region (900-1800 cm$^{-1}$

745        ) usually contains 235 wavenumbers. In that case, a window size of 5 points should

746        be used. The polynomial order for Savitzky-Golay fitting should be 2$^{nd}$ order for IR

747        spectroscopy due to the band shape.

748    **3. Light scattering correction using either multiplicative scatter correction (MSC),**

749        **SNV or 2$^{nd}$ derivative.** The user should prioritize MSC or SNV as these methods

750        maintain the spectral scale and original spectral shape. In case of unsatisfactory

751        results, 2$^{nd}$ derivative should be then employed.

752    **4. Baseline correction using automatic weighted least squares or rubber band**

753        **baseline correction.** If spectral differentiation is applied as light scattering correction

754        method, baseline correction is not necessary.

755    **5. Normalization** to the amide I peak, amide II peak or vector normalization (2-Norm,

756        length = 1) should be applied to correct different scales across spectra (*e.g.*, due to

757        different sample thicknesses when using FTIR in transmission mode).

758    **6. Scaling (*i.e.*, for each variable, mean-centring followed by division by the**

759        **variable standard deviation).** In case of data fusion, block-scaling should be used.

760 Data analysis

761 **(A) Exploratory analysis. ● TIMING 1h – 4 d (depending on the data size)**

762 Exploratory analysis should be primarily conducted using PCA. The PCA scores plot (PC1 vs

763 PC2) should be used for identification of the need of a standardization procedure.

764 **(B) Outlier detection. ● TIMING 1h – 1 d (depending on the data size)**

765 Apply PCA to the dataset and then estimate the Q residuals and Hotelling $T^2$ values. Use the

766 chart of Q residuals *versus* Hotelling $T^2$ to identify outliers. The outliers (*e.g.*, cosmic rays,

767 artefacts, low signal spectra and substrate only (non-tissue) spectra) should be removed from

768 the data set before proceeding to the next steps.

769 **(C) Sample split. ● TIMING 1 – 4 h (depending on the data size)**

770 Sample split should be performed before construction of standardization of multivariate

771 classification models. The samples can be split into training (70%) and test (30%) sets, using

772 a cross-validated model; or split into training (70%), validation (15%) and test (15%) sets

773 without using cross-validation. To maintain consistency and account for a well-balanced

774 training model, KS algorithm should be employed.

775 **(D) Standardization. ● TIMING 1h – 4 d (depending on the data size)**

776 **▲ CRITICAL** Standardization methods should be employed in the following order: DS >

777 PDS. The data from the secondary response should be separated into training (70%),

778 validation (15%) and test (15%) sets using KS algorithm. The number of transfer samples

779 should be firstly optimized using the validation set from the secondary response. Then, when

780 employing PDS, the window size should be optimized according to the size of the dataset.

781     (i) DS should be employed varying the number of transfer samples from 10-100% of

782 the training set from the primary system. The validation set from the secondary instrument

783    should be used to find the optimum number of transfer samples using the misclassification

784    rate as cost function.

785    (ii) PDS should be employed using the optimum number of samples found with DS.

786    Different window sizes should be tested using the validation set from the secondary system

787    with the misclassification rate as cost function. The window size should vary from 3-29 for a

788    spectral set with resolution of 4 cm$^{-1}$ in the biofingerprint region (235 variables).

789    **(E) Model construction. ● TIMING 1h – 4 d (depending on the data size)**

790    **▲ CRITICAL** Feature extraction (*e.g.*, by means of PCA) or feature selection (*e.g.*, by

791    means of GA or SPA) should be employed to reduce data collinearity and speed up data

792    processing and analysis time. PLS-DA is already a feature extraction method, thus the

793    performance of prior feature extraction is not necessary in this case. The classification

794    technique    employed    must    follow    a    parsimony    order:    LDA>PLS-

795    DA>QDA>KNN>SVM>ANN>Random forests>Deep learning approaches.

796    (i) Apply the feature extraction or selection technique. The optimization of the

797    number of PCs during PCA can be performed using an external validation set (15% of the

798    original data set) or using cross-validation (leave-one-out for small dataset [≤20 samples] or

799    venetian blinds [sample splitting: 10] for large datasets [>20 samples]). GA should be

800    realized three-times starting from different initial populations and the best result using an

801    external validation set (15% of the original data set) should be used. Cross-over probability

802    should be set for 40% and mutation probability should be set for 1-10% according to the size

803    of the dataset.

804    (ii) The classification method should be employed using optimization with an external

805    validation set or cross-validation, especially for selecting the number of latent variables of

806    PLS-DA and the kernel parameters for SVM. The kernel function for SVM should be RBF

807     kernel, due to its adaptation to different data distributions. To avoid overfitting, cross-

808     validation should be always performed during model construction to estimate the best RBF

809     parameters.

810     **? TROUBLESHOOTING**

811     **Spectral acquisition:** Spectral resolution, spectral range, SNR and signal aperture should be

812     optimized during experimental setup. Operators using different systems should try to keep

813     these parameters constant to reduce spectral differences.

814     **Data pre-processing:** To reduce spectral differences, the same data pre-processing should be

815     applied for spectra acquired in different systems.

816     **Standardization**: To improve the prediction capability of the classification model, the

817     primary system used should be the one with highest spectral resolution and smallest noise,

818     since all data from the secondary systems will be standardized to this pattern.

819     **● TIMING**

820     **Sample preparation:**

821     **(A)** Tissue (FFPE): 1-1.5 h

822     **(B)** Tissue (Snap-frozen or fresh): 2 h + drying time (3 h)

823     **(C)** Cells (fixed or live): 30 min + desiccation time (3 h)

824     **(D)** Biofluids (frozen or fresh): 5 min + thawing (20 min) + drying (1-1.5 h)

825     **Spectral acquisition:**   1 s – 5 min per spectrum (depending on the instrument and spectral
826     acquisition configurations)

827     **Data pre-processing:** 15 min – 4 h

828     **Data analysis:**

829     **(A)** Exploratory analysis: 1 h – 4 d

830     **(B)** Outlier detection: 1 h – 1 d

831 **(C)** Standardization: 1 h – 4 d

832 **(D)** Model construction: 1 h – 4 d

833 ANTICIPATED RESULTS

834       A pilot study was conducted to evaluate the effect of different instrument

835 manufacturers and operators towards spectral acquisition of healthy controls and ovarian

836 cancer samples based on blood plasma (5 healthy controls with 10 spectra per sample; 5

837 ovarian cancers with 10 spectra per sample) for a binary classification model using ATR-

838 FTIR spectroscopy. All specimens were collected with ethical approval obtained at Royal

839 Preston Hospital UK (16/EE/0010). Table 4 summarizes the experimental conditions in

840 which the experiments were performed.

841 **Table 4.** Experimental conditions for pilot study.

| Instrument | Operator | Spectral range | Number of co-additions | Spectral resolution | Room temperature | Air humidity |
|---|---|---|---|---|---|---|
| A | 1 | 4000-400 cm$^{-1}$ | 32 | 4 cm$^{-1}$ | 23.0ºC | 23% |
| | 2 | 4000-400 cm$^{-1}$ | 32 | 4 cm$^{-1}$ | 23.4ºC | 26% |
| B | 1 | 4000-400 cm$^{-1}$ | 32 | 4 cm$^{-1}$ | 24.0ºC | 26% |
| | 2 | 4000-400 cm$^{-1}$ | 32 | 4 cm$^{-1}$ | 24.9ºC | 24% |
| C | 1 | 4000-400 cm$^{-1}$ | 48 | 4 cm$^{-1}$ | 22.5ºC | 28% |
| | 2 | 4000-400 cm$^{-1}$ | 48 | 1 cm$^{-1}$ | 22.8ºC | 26% |

842

843       Instrument A and B were Bruker Tensor 27 with an HELIOS ATR attachment while

844 instrument C was an ATR-FTIR Thermo Scientific Nicolet iS10. The spectra were collected

845 for the same types of samples within three different days (operator 1: instrument A in day 1,

846 instrument B in day 3, and instrument C in day 2; operator 2: instrument A in day 2,

847 instrument B in day 1, and instrument C in day 3) and across two different laboratories

848 (instrument A and B in laboratory 1 and instrument C in laboratory 2). Each operator

849 prepared the samples individually from the same bulk, and measured them individually.

850 Spectral acquisition times were around 30 s for instruments A and B, and 40 s for instrument

851 C.

852     **(A) Effect of different instruments**

853         Three different ATR-FTIR spectrometers were used to analyse the samples. Data

854     were pre-processed by truncating at the biological fingerprint region (900-1800 cm$^{-1}$),

855     followed by Savitzky-Golay smoothing (window of 15 points, 2$^{nd}$ order polynomial

856     function), MSC, baseline correction using automatic weighted least squares and vector

857     normalization (2-Norm, length = 1). Each data set (A, B and C) was pre-processed

858     individually. The raw and pre-processed spectra for healthy controls and ovarian cancer

859     samples are depicted in Supplementary Material 1. All spectra collected by the three

860     instrument maintained the same spectral shape, indicating that the chemical information

861     stayed the same; however, large differences between the absorbance intensity were observed

862     between instrument C and the others (A, B), being caused due to different pressures applied

863     on the sample in the ATR module. The pressure applied to keep the sample in contact with

864     the ATR crystal directly affects the spectral signal intensity, which for instrument A and B

865     (same manufactures) were somewhere controlled by a contra weight, while for instrument C

866     the pressure was set based on a mechanical screw on the device, thus being biased by the

867     operator usage. The absorbance intensity variation between A and B is observed for this same

868     reason, but in a minor scale. Outlier detection was performed using a Hotelling T$^2$ versus Q

869     residual test (Supplementary Material 1).

870         **(i) Classification.** Classification was performed using PCA-LDA (10 PCs, explained

871     variance of 99.21%). Fig. 5a depicts the discriminant function (DF) score plot for PCA-LDA

872     using only the primary system (ATR-FTIR A). As observed, there is an almost perfect

873     separation between the samples from the two classes (accuracy = 100%, sensitivity = 100%,

874     specificity = 100%). However, when the spectra acquired using instruments B and C are

875     predicted using the model for A, the results decreased significantly (accuracy = 66.7%,

876  sensitivity = 83.2%, specificity = 48.9%) (Fig. 5b), necessitating the use of a standardization
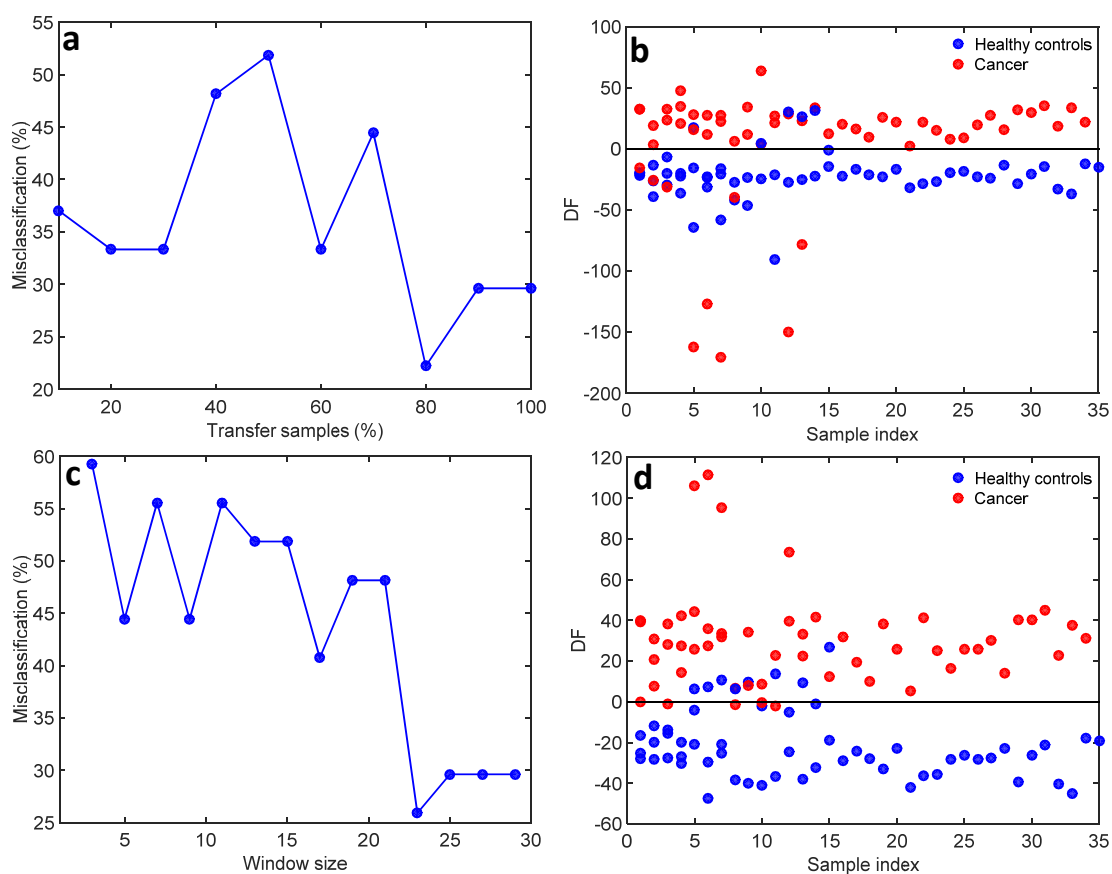
877  procedure.



878

879  **Figure 5.** (a) DF plot of the PCA-LDA model for the primary system; (b) DF plot of the

880  PCA-LDA model for the primary system predicting the samples from the secondary systems.

881

882      **(ii) Standardization.** Standardization was employed using both DS and PDS in order

883  to compare the two methods. The number of transfer samples for DS was optimized

884  according to the misclassification rate obtained for the validation set using the secondary

885  system (Fig. 6a). An optimum number corresponding to 80% of the samples in the training

886  set of the primary system (55 transfer samples) was obtained, resulting to a misclassification

887  rate of 22.2% in the validation set of the secondary system. This improved the accuracy

888  (77.8%) and specificity (80.0%). Sensitivity decreased to 75.0%, which is an acceptable

889  value. The results after DS are better balanced than without standardization. Fig. 6b shows

890  the DF plot for the PCA-LDA model using the training of the primary system and prediction

891  with the secondary system after DS.

892      PDS was also applied. The number of transfer samples was maintained as 55 (80% of

893  the primary training set) and the window size was optimized by using the validation set of the

45

secondary system. An optimum window size of 23 wavenumbers was selected with a misclassification rate of 25.9% (Fig. 6c). The accuracy, sensitivity and specificity using PDS were 74.1%, 71.4% and 75.0%, respectively. The DS presented a slightly higher performance than PDS for this dataset. However, DS generated some outliers not observed before, while PDS did not. Thus, in general, PDS provided a better standardization of the data. The PCA-LDA DF plot after PDS is depicted in Fig. 6d.



**Figure 6.** (a) Misclassification rate in % for the validation set of the secondary system varying the number of transfer samples in % from the primary system for DS optimization; (b) DF plot of the PCA-LDA model for the primary system predicting the validation set from the secondary system after DS; (c) Misclassification rate in % for the validation set of the secondary system varying the window size for PDS optimization; (d) DF plot of the PCA-

906 LDA model for the primary system predicting the validation set from the secondary system

907 after PDS.

908 **(B) Effect of different operators**

909       The effect of different user operators acquiring spectra from the same samples using

910 the same instruments was also evaluated. Similarly to before, data were pre-processed by

911 cutting the biological fingerprint region (900-1800 cm$^{-1}$), followed by Savitzky-Golay

912 smoothing (window of 15 points, 2$^{nd}$ order polynomial function), MSC, baseline correction

913 using automatic weighted least squares and vector normalization (2-Norm, length = 1). Each

914 dataset was pre-processed individually. All raw and pre-processed spectra varying operators

915 are depicted in Supplementary Material 1. Outlier detection was performed using a Hotelling

916 T$^2$ versus Q residual test (Supplementary Material 1). The PCA scores plots for the pre-

917 processed spectra are depicted in Supplementary Material 1. The main difference between the

918 operators was observed for instrument C (Supplementary Material 1, Figure S5e), since the

919 spectral resolutions used by them were different, which can cause major data distortion.

920       **(i) Classification.** Classification was performed using PCA-LDA (10 PCs, explained

921 variance of 98.62%). Fig. 7a depicts the DF score plot for PCA-LDA using only the primary

922 system (Operator 1). There is a significant separation between the samples from the two

923 classes (accuracy = 88.4%, sensitivity = 77.3%, specificity = 100%). When the spectra

924 acquired by Operator 2 are predicted using the model for Operator 1, the results decreased

925 (accuracy = 75.6%, sensitivity = 66.7%, specificity = 84.6%) (Fig. 7b), which again

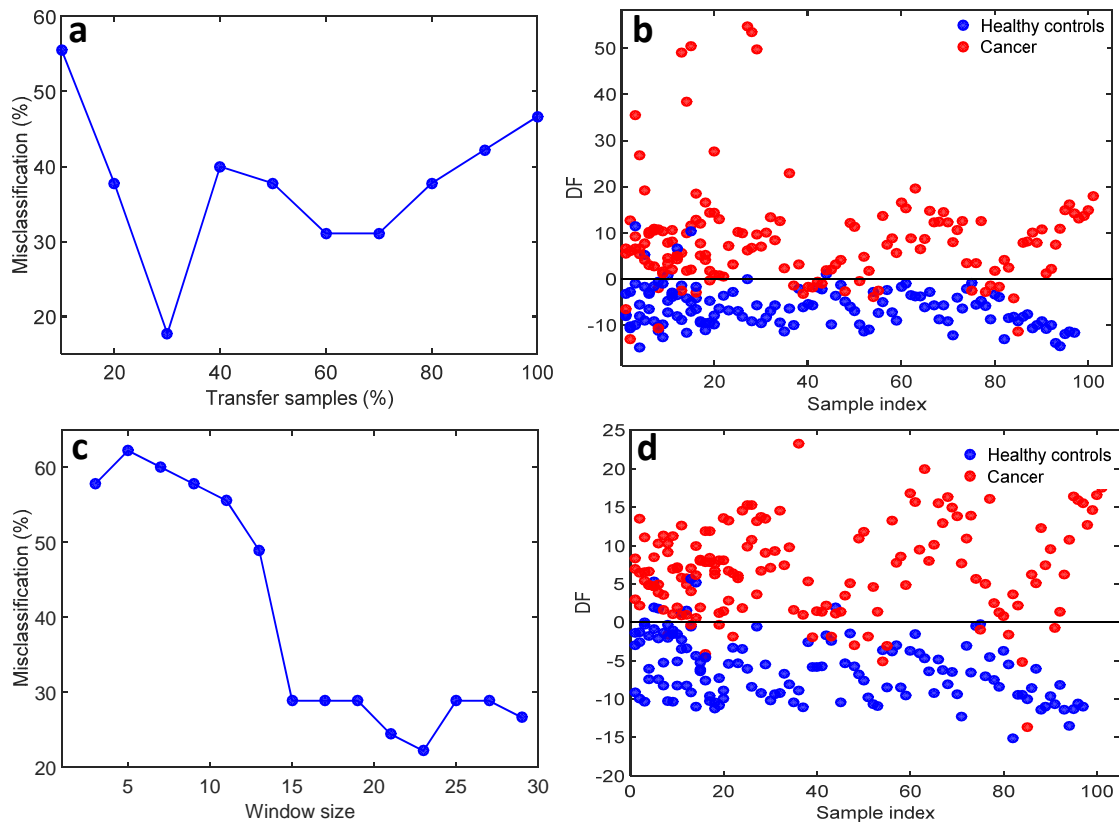926 necessitates the use of a standardization procedure.

927

928

929

**Figure 7.** (a) DF plot of the PCA-LDA model for the primary system (Operator 1); (b) DF plot of the PCA-LDA model for the primary system predicting the samples from the secondary system (Operator 2).

933

**(ii) Standardization.** DS and PDS were employed as standardization methods. The number of transfer samples for DS was optimized according to the misclassification rate obtained for the validation set using the secondary system (Operator 2) (Fig. 8a). An optimum number of 59 transfer samples (30% of the samples in the training set of the primary system [Operator 1]) was obtained, resulting in a misclassification rate of 17.8% in the validation set of the secondary system. This improved the accuracy (82.2%), sensitivity (69.6%) and specificity (95.5%) compared to the results without DS. Fig. 8b shows the DF plot for the PCA-LDA model using the training of the primary system and prediction with the secondary system after DS.

The number of transfer samples was maintained as 59 for PDS; and the window size was optimized by using the validation set of the secondary system. An optimum window size of 23 wavenumbers was selected with a misclassification rate of 22.2% (Fig. 8c). The

48

946    accuracy, sensitivity and specificity using PDS were 77.8%, 100% and 54.5%, respectively.

947    Although DS obtained an average better classification performance than PDS for this dataset,

948    it also generated some outliers as mentioned before. For this reason, the results after PDS

949    seem better standardized. The PCA-LDA DF plot after PDS is depicted in Fig. 8d.



950

951    **Figure 8.** (a) Misclassification rate in % for the validation set of the secondary system

952    (Operator 2) varying the number of transfer samples in % from the primary system (Operator

953    1) for DS optimization; (b) DF plot of the PCA-LDA model for the primary system

954    predicting the validation set from the secondary system after DS; (c) Misclassification rate in

955    % for the validation set of the secondary system varying the window size for PDS

956    optimization; (d) DF plot of the PCA-LDA model for the primary system predicting the

957    validation set from the secondary system after PDS.

958

## Acknowledgements

## Author contributions

F.L.M. is the principal investigator who conceived the idea for the manuscript; C.L.M.M. and M.P. wrote the manuscript. All co-authors contributed recommendations and provided feedback and changes to the manuscript; and, C.L.M.M., M.P. and F.L.M. brought together the text and finalized the manuscript.

## Competing financial interests

The authors declare no competing financial interest.

## Data availability statement

The datasets generated during and/or analysed during the current study are available from the corresponding author on reasonable request.

References

| 974 | References |
| --- | --- |

1    Hofmann-Wellenhof, B., Lichtenegger, H. & Collins, J. *Global positioning system: theory and practice*.  (Springer Science & Business Media, 2012).

2    Morris, P. & Perkins, A. Diagnostic imaging. *Lancet* **379**, 1525-1533 (2012).

3    Lee, S. S. *et al.* Crohn disease of the small bowel: comparison of CT enterography, MR enterography, and small-bowel follow-through as diagnostic techniques. *Radiology* **251**, 751-761 (2009).

4    Lagleyre, S. *et al.* Reliability of high-resolution CT scan in diagnosis of otosclerosis. *Otol Neurotol* **30**, 1152-1159 (2009).

5    Kalita, J. & Misra, U. Comparison of CT scan and MRI findings in the diagnosis of Japanese encephalitis. *J Neurol Sci* **174**, 3-8 (2000).

6    Schrevens, L., Lorent, N., Dooms, C. & Vansteenkiste, J. The role of PET scan in diagnosis, staging, and management of non-small cell lung cancer. *Oncologist* **9**, 633-643 (2004).

7    Jagust, W., Reed, B., Mungas, D., Ellis, W. & Decarli, C. What does fluorodeoxyglucose PET imaging add to a clinical diagnosis of dementia? *Neurology* **69**, 871-877 (2007).

8    Zhou, M. *et al.* Clinical utility of breast-specific gamma imaging for evaluating disease extent in the newly diagnosed breast cancer patient. *Am J Surg* **197**, 159-163 (2009).

9    Wallace, B. A. *et al.* Biomedical applications of synchrotron radiation circular dichroism spectroscopy: identification of mutant proteins associated with disease and development of a reference database for fold motifs. *Faraday Discuss* **126**, 237-243 (2004).

10   Greenfield, N. J. Using circular dichroism spectra to estimate protein secondary structure. *Nat Protoc* **1**, 2876 (2006).

11   Micsonai, A. *et al.* Accurate secondary structure prediction and fold recognition for circular dichroism spectroscopy. *Proc Natl Acad Sci USA* **112**, E3095-E3103 (2015).

12   Miles, A. J. & Wallace, B. A. Circular dichroism spectroscopy of membrane proteins. *Chem Soc Rev* **45**, 4859-4872 (2016).

13   Brown, J. Q., Vishwanath, K., Palmer, G. M. & Ramanujam, N. Advances in quantitative UV–visible spectroscopy for clinical and pre-clinical application in cancer. *Curr Opin Biotechnol* **20**, 119-131 (2009).

14   Yang, P.-W. *et al.* Visible-absorption spectroscopy as a biomarker to predict treatment response and prognosis of surgically resected esophageal cancer. *Sci Rep* **6**, 33414 (2016).

15   Organization, W. H. *Fluorescence microscopy for disease diagnosis and environmental monitoring*.  (2005).

16   Shahzad, A. *et al.* Diagnostic application of fluorescence spectroscopy in oncology field: hopes and challenges. *Appl Spectrosc Rev* **45**, 92-99 (2010).

17   Sieroń, A. *et al.* The role of fluorescence diagnosis in clinical practice. *Onco Targets Ther* **6**, 977 (2013).

18   Shin, D., Vigneswaran, N., Gillenwater, A. & Richards-Kortum, R. Advances in fluorescence imaging techniques to detect oral cancer and its precursors. *Future Oncol* **6**, 1143-1154 (2010).

19   Shahzad, A. *et al.* Emerging applications of fluorescence spectroscopy in medical microbiology field. *J Transl Med* **7**, 99 (2009).

20   Möller-Hartmann, W. *et al.* Clinical application of proton magnetic resonance spectroscopy in the diagnosis of intracranial mass lesions. *Neuroradiology* **44**, 371-381 (2002).

21   Gowda, G. N. *et al.* Metabolomics-based methods for early disease diagnostics. *Expert Rev Mol Diagn* **8**, 617-633 (2008).

22   Frisoni, G. B., Fox, N. C., Jack, C. R., Scheltens, P. & Thompson, P. M. The clinical use of structural MRI in Alzheimer disease. *Nat Rev Neurol* **6**, 67-77 (2010).

23   Chan, A. W. *et al.* 1 H-NMR urinary metabolomic profiling for diagnosis of gastric cancer. *Br J Cancer* **114**, 59 (2016).

1024  24  Palmnas, M. S. & Vogel, H. J. The future of NMR metabolomics in cancer therapy: towards
1025      personalizing treatment and developing targeted drugs? *Metabolites* **3**, 373-396 (2013).
1026  25  Patil, P. & Dasgupta, B. Role of diagnostic ultrasound in the assessment of musculoskeletal
1027      diseases. *Ther Adv Musculoskelet Dis* **4**, 341-355 (2012).
1028  26  Navani, N. *et al.* Lung cancer diagnosis and staging with endobronchial ultrasound-guided
1029      transbronchial needle aspiration compared with conventional approaches: an open-label,
1030      pragmatic, randomised controlled trial. *Lancet Respir Med* **3**, 282-289 (2015).
1031  27  Menon, U. *et al.* Sensitivity and specificity of multimodal and ultrasound screening for
1032      ovarian cancer, and stage distribution of detected cancers: results of the prevalence screen
1033      of the UK Collaborative Trial of Ovarian Cancer Screening (UKCTOCS). *Lancet Oncol* **10**, 327-
1034      340 (2009).
1035  28  Smith-Bindman, R. *et al.* Endovaginal ultrasound to exclude endometrial cancer and other
1036      endometrial abnormalities. *Jama* **280**, 1510-1517 (1998).
1037  29  Gajjar, K. *et al.* Diagnostic segregation of human brain tumours using Fourier-transform
1038      infrared and/or Raman spectroscopy coupled with discriminant analysis. *Anal Methods* **5**,
1039      89-102 (2013).
1040  30  Bury, D. *et al.* Phenotyping Metastatic Brain Tumors Applying Spectrochemical Analyses:
1041      Segregation of Different Cancer Types. *Anal. Lett.*, 1-2 (2018).
1042  31  Hands, J. R. *et al.* Attenuated Total Reflection Fourier Transform Infrared (ATR-FTIR) spectral
1043      discrimination of brain tumour severity from serum samples. *J Biophotonics* **7**, 189-199
1044      (2014).
1045  32  Hands, J. R. *et al.* Brain tumour differentiation: rapid stratified serum diagnostics via
1046      attenuated total reflection Fourier-transform infrared spectroscopy. *Journal of neuro-
1047      oncology* **127**, 463-472 (2016).
1048  33  Walsh, M. J., Kajdacsy-Balla, A., Holton, S. E. & Bhargava, R. Attenuated total reflectance
1049      Fourier-transform infrared spectroscopic imaging for breast histopathology. *Vib Spectrosc*
1050      **60**, 23-28 (2012).
1051  34  Lane, R. & Seo, S. S. Attenuated Total Reflectance Fourier Transform Infrared Spectroscopy
1052      Method to Differentiate Between Normal and Cancerous Breast Cells. *J Nanosci Nanotechnol*
1053      **12**, 7395-7400 (2012).
1054  35  Backhaus, J. *et al.* Diagnosis of breast cancer with infrared spectroscopy from serum
1055      samples. *Vib Spectrosc* **52**, 173-177 (2010).
1056  36  Wang, J.-S. *et al.* FT-IR spectroscopic analysis of normal and cancerous tissues of esophagus.
1057      *World journal of gastroenterology* **9**, 1897 (2003).
1058  37  Maziak, D. E. *et al.* Fourier-transform infrared spectroscopic study of characteristic
1059      molecular structure in cancer cells of esophagus: an exploratory study. *Cancer Detect. Prev.*
1060      **31** (2007).
1061  38  McIntosh, L. M. *et al.* Infrared spectra of basal cell carcinomas are distinct from non-tumor-
1062      bearing skin components. *J Investig Dermatol* **112**, 951-956 (1999).
1063  39  McIntosh, L. M. *et al.* Towards non-invasive screening of skin lesions by near-infrared
1064      spectroscopy. *Journal of Investigative Dermatology* **116**, 175-181 (2001).
1065  40  Mostaço-Guidolin, L. B., Murakami, L. S., Nomizo, A. & Bachmann, L. Fourier transform
1066      infrared spectroscopy of skin cancer cells and tissues. *Appl Spectrosc Rev* **44**, 438-455 (2009).
1067  41  Mordechai, S. *et al.* Possible common biomarkers from FTIR microspectroscopy of cervical
1068      cancer and melanoma. *Journal of microscopy* **215**, 86-91 (2004).
1069  42  Hammody, Z., Sahu, R. K., Mordechai, S., Cagnano, E. & Argov, S. Characterization of
1070      malignant melanoma using vibrational spectroscopy. *The Scientific World Journal* **5**, 173-182
1071      (2005).
1072  43  Kondepati, V. R., Keese, M., Mueller, R., Manegold, B. C. & Backhaus, J. Application of near-
1073      infrared spectroscopy for the diagnosis of colorectal cancer in resected human tissue
1074      specimens. *Vib Spectrosc* **44**, 236-242 (2007).

1075    44    Rigas, B., Morgello, S., Goldman, I. S. & Wong, P. Human colorectal cancers display abnormal
1076        Fourier-transform infrared spectra. *Proceedings of the National Academy of Sciences* **87**,
1077        8140-8144 (1990).
1078    45    Yao, H., Shi, X. & Zhang, Y. The Use of FTIR-ATR Spectrometry for Evaluation of Surgical
1079        Resection Margin in Colorectal Cancer: A Pilot Study of 56 Samples. *J Spectrosc* **2014**, 4
1080        (2014).
1081    46    Lewis, P. D. *et al.* Evaluation of FTIR Spectroscopy as a diagnostic tool for lung cancer using
1082        sputum. *BMC Cancer* **10**, 640 (2010).
1083    47    Akalin, A. *et al.* Classification of malignant and benign tumors of the lung by infrared spectral
1084        histopathology (SHP). *Lab Invest* **95**, 406 (2015).
1085    48    Großerueschkamp, F. *et al.* Marker-free automated histopathological annotation of lung
1086        tumour subtypes by FTIR imaging. *Analyst* **140**, 2114-2120 (2015).
1087    49    Owens, G. L. *et al.* Vibrational biospectroscopy coupled with multivariate analysis extracts
1088        potentially diagnostic features in blood plasma/serum of ovarian cancer patients. *J*
1089        *Biophotonics* **7**, 200-209 (2014).
1090    50    Gajjar, K. *et al.* Fourier-transform infrared spectroscopy coupled with a classification
1091        machine for the analysis of blood plasma or serum: a novel diagnostic approach for ovarian
1092        cancer. *Analyst* **138**, 3917-3926 (2013).
1093    51    Theophilou, G., Lima, K. M. G., Martin-Hirsch, P. L., Stringfellow, H. F. & Martin, F. L. ATR-
1094        FTIR spectroscopy coupled with chemometric analysis discriminates normal, borderline and
1095        malignant ovarian tissue: classifying subtypes of human cancer. *Analyst* **141**, 585-594 (2016).
1096    52    Mehrotra, R., Tyagi, G., Jangir, D. K., Dawar, R. & Gupta, N. Analysis of ovarian tumor
1097        pathology by Fourier Transform Infrared Spectroscopy. *J Ovarian Res* **3**, 27 (2010).
1098    53    Paraskevaidi, M. *et al.* Potential of mid-infrared spectroscopy as a non-invasive diagnostic
1099        test in urine for endometrial or ovarian cancer. *Analyst* (2018).
1100    54    Taylor, S. E. *et al.* Infrared spectroscopy with multivariate analysis to interrogate
1101        endometrial tissue: a novel and objective diagnostic approach. *Br J Cancer* **104**, 790-797
1102        (2011).
1103    55    Paraskevaidi, M. *et al.* Aluminium foil as an alternative substrate for the spectroscopic
1104        interrogation of endometrial cancer. *J Biophotonics* (2018).
1105    56    Gajjar, K. *et al.* Histology verification demonstrates that biospectroscopy analysis of cervical
1106        cytology identifies underlying disease more accurately than conventional screening:
1107        removing the confounder of discordance. *PLoS One* **9**, e82416 (2014).
1108    57    Walsh, M. J. *et al.* IR microspectroscopy: potential applications in cervical cancer screening.
1109        *Cancer Lett.* **246**, 1-11 (2007).
1110    58    Wood, B. R., Quinn, M. A., Burden, F. R. & McNaughton, D. An investigation into FTIR
1111        spectroscopy as a biodiagnostic tool for cervical cancer. *Biospectroscopy* **2**, 143-153 (1996).
1112    59    Podshyvalov, A. *et al.* Distinction of cervical cancer biopsies by use of infrared
1113        microspectroscopy and probabilistic neural networks. *Appl Opt* **44**, 3725-3734 (2005).
1114    60    Theophilou, G. *et al.* A biospectroscopic analysis of human prostate tissue obtained from
1115        different time periods points to a trans-generational alteration in spectral phenotype. *Sci*
1116        *Rep* **5**, 13465 (2015).
1117    61    Baker, M. J. *et al.* Investigating FTIR based histopathology for the diagnosis of prostate
1118        cancer. *J Biophotonics* **2** (2009).
1119    62    Derenne, A., Gasper, R. & Goormaghtigh, E. The FTIR spectrum of prostate cancer cells
1120        allows the classification of anticancer drugs according to their mode of action. *Analyst* **136**
1121        (2011).
1122    63    Gazi, E. *et al.* A correlation of FTIR spectra derived from prostate cancer biopsies with
1123        Gleason grade and tumour stage. *European urology* **50**, 750-761 (2006).
1124    64    Paraskevaidi, M. *et al.* Differential diagnosis of Alzheimer's disease using spectrochemical
1125        analysis of blood. *Proc Natl Acad Sci USA*, 201701517 (2017).

| 1126 | 65 | Carmona, P. *et al.* Discrimination analysis of blood plasma associated with Alzheimer's |
|------|----|---|
| 1127 | | disease using vibrational spectroscopy. *J Alzheimers Dis* **34**, 911-920 (2013). |
| 1128 | 66 | Carmona, P., Molina, M., López-Tobar, E. & Toledano, A. Vibrational spectroscopic analysis |
| 1129 | | of peripheral blood plasma of patients with Alzheimer's disease. *Anal Bioanal Chem* **407**, |
| 1130 | | 7747-7756 (2015). |
| 1131 | 67 | Paraskevaidi, M. *et al.* Blood-based near-infrared spectroscopy for the rapid low-cost |
| 1132 | | detection of Alzheimer's disease. *Analyst* (2018). |
| 1133 | 68 | Sitole, L., Steffens, F., Krüger, T. P. J. & Meyer, D. Mid-ATR-FTIR Spectroscopic Profiling of |
| 1134 | | HIV/AIDS Sera for Novel Systems Diagnostics in Global Health. *OMICS* **18**, 513-523 (2014). |
| 1135 | 69 | Coopman, R. *et al.* Glycation in human fingernail clippings using ATR-FTIR spectrometry, a |
| 1136 | | new marker for the diagnosis and monitoring of diabetes mellitus. *Clin Biochem* **50**, 62-67 |
| 1137 | | (2017). |
| 1138 | 70 | Scott, D. A. *et al.* Diabetes-related molecular signatures in infrared spectra of human saliva. |
| 1139 | | *Diabetol Metab Syndr* **2**, 48 (2010). |
| 1140 | 71 | Varma, V. K., Kajdacsy-Balla, A., Akkina, S. K., Setty, S. & Walsh, M. J. A label-free approach |
| 1141 | | by infrared spectroscopic imaging for interrogating the biochemistry of diabetic |
| 1142 | | nephropathy progression. *Kidney Int* **89**, 1153-1159 (2016). |
| 1143 | 72 | Lechowicz, L., Chrapek, M., Gaweda, J., Urbaniak, M. & Konieczna, I. Use of Fourier- |
| 1144 | | transform infrared spectroscopy in the diagnosis of rheumatoid arthritis: a pilot study. *Mol* |
| 1145 | | *Biol Rep* **43**, 1321-1326 (2016). |
| 1146 | 73 | Canvin, J. *et al.* Infrared spectroscopy: shedding light on synovitis in patients with |
| 1147 | | rheumatoid arthritis. *Rheumatology* **42**, 76-82 (2003). |
| 1148 | 74 | Oemrawsingh, R. M. *et al.* Near-infrared spectroscopy predicts cardiovascular outcome in |
| 1149 | | patients with coronary artery disease. *J Am Coll Cardiol* **64**, 2510-2518 (2014). |
| 1150 | 75 | Wang, J. *et al.* Near-infrared spectroscopic characterization of human advanced |
| 1151 | | atherosclerotic plaques. *J Am Coll Cardiol* **39**, 1305-1313 (2002). |
| 1152 | 76 | Martin, M. *et al.* The effect of common anticoagulants in detection and quantification of |
| 1153 | | malaria parasitemia in human red blood cells by ATR-FTIR spectroscopy. *Analyst* (2017). |
| 1154 | 77 | Khoshmanesh, A. *et al.* Detection and Quantification of Early-Stage Malaria Parasites in |
| 1155 | | Laboratory Infected Erythrocytes by Attenuated Total Reflectance Infrared Spectroscopy and |
| 1156 | | Multivariate Analysis. *Anal Chem* **86**, 4379-4386 (2014). |
| 1157 | 78 | Roy, S. *et al.* Simultaneous ATR-FTIR Based Determination of Malaria Parasitemia, Glucose |
| 1158 | | and Urea in Whole Blood Dried onto a Glass Slide. *Anal Chem* **89**, 5238-5245 (2017). |
| 1159 | 79 | Markus, A. P. J. *et al.* New technique for diagnosis and monitoring of alcaptonuria: |
| 1160 | | quantification of homogentisic acid in urine with mid-infrared spectrometry. *Anal Chim Acta* |
| 1161 | | **429**, 287-292 (2001). |
| 1162 | 80 | Grimard, V. *et al.* Phosphorylation-induced Conformational Changes of Cystic Fibrosis |
| 1163 | | Transmembrane Conductance Regulator Monitored by Attenuated Total Reflection-Fourier |
| 1164 | | Transform IR Spectroscopy and Fluorescence Spectroscopy. *J Biol Chem* **279**, 5528-5536 |
| 1165 | | (2004). |
| 1166 | 81 | Aksoy, C., Guliyev, A., Kilic, E., Uckan, D. & Severcan, F. Bone marrow mesenchymal stem |
| 1167 | | cells in patients with beta thalassemia major: molecular analysis with attenuated total |
| 1168 | | reflection-Fourier transform infrared spectroscopy study as a novel method. *Stem Cells Dev* |
| 1169 | | **21**, 2000-2011 (2012). |
| 1170 | 82 | Graça, G. *et al.* Mid-infrared (MIR) metabolic fingerprinting of amniotic fluid: A possible |
| 1171 | | avenue for early diagnosis of prenatal disorders? *Anal Chim Acta* **764**, 24-31 (2013). |
| 1172 | 83 | Hasegawa, J. *et al.* Evaluation of placental function using near infrared spectroscopy during |
| 1173 | | fetal growth restriction. *J Perinatal Med* **38**, 29-32 (2010). |
| 1174 | 84 | Theelen, T., Berendschot, T. T., Hoyng, C. B., Boon, C. J. & Klevering, B. J. Near-infrared |
| 1175 | | reflectance imaging of neovascular age-related macular degeneration. *Graefe's Archive for* |
| 1176 | | *Clinical and Experimental Ophthalmology* **247**, 1625 (2009). |

1177 85 Semoun, O. *et al.* Infrared features of classic choroidal neovascularisation in exudative age-
1178     related macular degeneration. *Br. J. Ophthalmol.* **93**, 182-185 (2009).
1179 86 Peters, A. S. *et al.* Serum-infrared spectroscopy is suitable for diagnosis of atherosclerosis
1180     and its clinical manifestations. *Vib Spectrosc* **92**, 20-26 (2017).
1181 87 Afara, I. O., Prasadam, I., Arabshahi, Z., Xiao, Y. & Oloyede, A. Monitoring osteoarthritis
1182     progression using near infrared (NIR) spectroscopy. *Sci Rep* **7**, 11463 (2017).
1183 88 Bi, X. *et al.* Fourier transform infrared imaging and MR microscopy studies detect
1184     compositional and structural changes in cartilage in a rabbit model of osteoarthritis. *Anal*
1185     *Bioanal Chem* **387**, 1601-1612 (2007).
1186 89 David-Vaudey, E. *et al.* Fourier Transform Infrared Imaging of focal lesions in human
1187     osteoarthritic cartilage. *Eur Cell Mater* **10**, 60 (2005).
1188 90 Trevisan, J., Angelov, P. P., Carmichael, P. L., Scott, A. D. & Martin, F. L. Extracting biological
1189     information with computational analysis of Fourier-transform infrared (FTIR)
1190     biospectroscopy datasets: current practices to future perspectives. *Analyst* **137**, 3202-3215
1191     (2012).
1192 91 Baker, M. J. *et al.* Using Fourier transform IR spectroscopy to analyze biological materials.
1193     *Nat Protoc* **9**, 1771-1791 (2014).
1194 92 Andrew Chan, K. L. & Kazarian, S. G. Attenuated total reflection Fourier-transform infrared
1195     (ATR-FTIR) imaging of tissues and live cells. *Chem Soc Rev* **45**, 1850-1864 (2016).
1196 93 Pilling, M. & Gardner, P. Fundamental developments in infrared spectroscopic imaging for
1197     biomedical applications. *Chem Soc Rev* **45**, 1935-1957 (2016).
1198 94 Martin, F. L. *et al.* Distinguishing cell types or populations based on the computational
1199     analysis of their infrared spectra. *Nat Protoc* **5**, 1748-1760 (2010).
1200 95 Butler, H. J. *et al.* Using Raman spectroscopy to characterize biological materials. *Nat Protoc*
1201     **11**, 664-687 (2016).
1202 96 Kong, L. *et al.* Characterization of bacterial spore germination using phase-contrast and
1203     fluorescence microscopy, Raman spectroscopy and optical tweezers. *Nat Protoc* **6**, 625
1204     (2011).
1205 97 Harmsen, S., Wall, M. A., Huang, R. & Kircher, M. F. Cancer imaging using surface-enhanced
1206     resonance Raman scattering nanoparticles. *Nat Protoc* **12**, 1400 (2017).
1207 98 Beckonert, O. *et al.* Metabolic profiling, metabolomic and metabonomic procedures for
1208     NMR spectroscopy of urine, plasma, serum and tissue extracts. *Nat Protoc* **2**, 2692 (2007).
1209 99 Felten, J. *et al.* Vibrational spectroscopic image analysis of biological material using
1210     multivariate curve resolution–alternating least squares (MCR-ALS). *Nat Protoc* **10**, 217
1211     (2015).
1212 100 Yang, H., Yang, S., Kong, J., Dong, A. & Yu, S. Obtaining information about protein secondary
1213     structures in aqueous solution using Fourier transform IR spectroscopy. *Nat Protoc* **10**, 382
1214     (2015).
1215 101 Sreedhar, H. *et al.* High-definition Fourier transform infrared (FT-IR) spectroscopic imaging of
1216     human tissue sections towards improving pathology. *J Vis Exp* (2015).
1217 102 Varriale, A. *et al.* Fluorescence correlation spectroscopy assay for gliadin in food. *Anal Chem*
1218     **79**, 4687-4689 (2007).
1219 103 Song, X., Li, H., Al-Qadiri, H. M. & Lin, M. Detection of herbicides in drinking water by
1220     surface-enhanced Raman spectroscopy coupled with gold nanostructures. *J Food Meas*
1221     *Charact* **7**, 107-113 (2013).
1222 104 Osborne, B. G. & Fearn, T. Near-infrared spectroscopy in food analysis. *Encyclopedia Anal*
1223     *Chem* **5**, 4069-4082 (2000).
1224 105 Qu, J.-H. *et al.* Applications of near-infrared spectroscopy in food safety evaluation and
1225     control: A review of recent research advances. *Crit. Rev. Food Sci. Nutr.* **55**, 1939-1954
1226     (2015).

1227 106 Penido, C. A. F., Pacheco, M. T. T., Lednev, I. K. & Silveira, L. Raman spectroscopy in forensic
1228    analysis: identification of cocaine and other illegal drugs of abuse. *J Raman Spectrosc* **47**, 28-
1229    38 (2016).
1230 107 Ryder, A. G. Classification of narcotics in solid mixtures using principal component analysis
1231    and Raman spectroscopy. *J Forensic Sci* **47**, 275-284 (2002).
1232 108 Melin, A. M., Perromat, A. & Déléris, G. Pharmacologic application of Fourier transform IR
1233    spectroscopy: in vivo toxicity of carbon tetrachloride on rat liver. *Biopolymers: Original*
1234    *Research on Biomolecules* **57**, 160-168 (2000).
1235 109 Harrigan, G. G. *et al.* Application of high-throughput Fourier-transform infrared spectroscopy
1236    in toxicology studies: contribution to a study on the development of an animal model for
1237    idiosyncratic toxicity. *Toxicol. Lett.* **146**, 197-205 (2004).
1238 110 Choo-Smith, L.-P. *et al.* Investigating microbial (micro) colony heterogeneity by vibrational
1239    spectroscopy. *Appl Environ Microbiol* **67**, 1461-1469 (2001).
1240 111 Helm, D., Labischinski, H., Schallehn, G. & Naumann, D. Classification and identification of
1241    bacteria by Fourier-transform infrared spectroscopy. *Microbiology* **137**, 69-79 (1991).
1242 112 Carmona, P., Monzon, M., Monleon, E., Badiola, J. J. & Monreal, J. In vivo detection of
1243    scrapie cases from blood by infrared spectroscopy. *J. Gen. Virol.* **86**, 3425-3431 (2005).
1244 113 Cui, L. *et al.* A novel functional single-cell approach to probing nitrogen-fixing bacteria in soil
1245    communities by resonance Raman spectroscopy with 15N2 labelling. *Anal. Chem.*
1246    **10.1021/acs.analchem.7b05080.** (2018).
1247 114 Lasch, P. & Naumann, D. Infrared spectroscopy in microbiology. *Encyclopedia Anal Chem*
1248    (2015).
1249 115 Maquelin, K. *et al.* Identification of medically relevant microorganisms by vibrational
1250    spectroscopy. *J Microbiol Methods* **51**, 255-271 (2002).
1251 116 Day, J. S., Edwards, H. G., Dobrowski, S. A. & Voice, A. M. The detection of drugs of abuse in
1252    fingerprints using Raman spectroscopy I: latent fingerprints. *Spectrochim Acta A Mol Biomol*
1253    *Spectrosc* **60**, 563-568 (2004).
1254 117 Macleod, N. A. & Matousek, P. Emerging Non-invasive Raman Methods in Process Control
1255    and Forensic Applications. *Pharm Res* **25**, 2205 (2008).
1256 118 Lewis, I., Daniel Jr, N., Chaffin, N., Griffiths, P. & Tungol, M. Raman spectroscopic studies of
1257    explosive materials: towards a fieldable explosives detector. *Spectrochimica Acta Part A:*
1258    *Molecular and Biomolecular Spectroscopy* **51**, 1985-2000 (1995).
1259 119 Hargreaves, M. D. & Matousek, P. Threat detection of liquid explosive precursor mixtures by
1260    Spatially Offset Raman Spectroscopy (SORS). in *Optics and photonics for counterterrorism*
1261    *and crime fighting V.* Vol. **7486** 74860B (International Society for Optics and Photonics).
1262 120 Ali, E. M., Edwards, H. G., Hargreaves, M. D. & Scowen, I. J. Raman spectroscopic
1263    investigation of cocaine hydrochloride on human nail in a forensic context. *Anal Bioanal*
1264    *Chem* **390**, 1159-1166 (2008).
1265 121 Vergote, G. J., Vervaet, C., Remon, J. P., Haemers, T. & Verpoort, F. Near-infrared FT-Raman
1266    spectroscopy as a rapid analytical tool for the determination of diltiazem hydrochloride in
1267    tablets. *Eur. J. Pharm. Sci.* **16**, 63-67 (2002).
1268 122 Eliasson, C. & Matousek, P. Noninvasive authentication of pharmaceutical products through
1269    packaging using spatially offset Raman spectroscopy. *Anal Chem* **79**, 1696-1701 (2007).
1270 123 Lohr, D. *et al.* Non-destructive determination of carbohydrate reserves in leaves of
1271    ornamental cuttings by near-infrared spectroscopy (NIRS) as a key indicator for quality
1272    assessments. *Biosys Eng* **158**, 51-63 (2017).
1273 124 Heys, K. A., Shore, R. F., Pereira, M. G. & Martin, F. L. Levels of Organochlorine Pesticides Are
1274    Associated with Amyloid Aggregation in Apex Avian Brains. *Environ Sci Technol* **51**, 8672-
1275    8681 (2017).

| 1276 | 125 | Comino, F., Aranda, V., García-Ruiz, R. & Domínguez-Vidal, A. Infrared spectroscopy as a tool |
| 1277 | | for the assessment of soil biological quality in agricultural soils under contrasting |
| 1278 | | management practices. *Ecol Indicators* **87**, 117-126 (2018). |
| 1279 | 126 | Eliasson, C., Macleod, N. & Matousek, P. Noninvasive detection of concealed liquid |
| 1280 | | explosives using Raman spectroscopy. *Anal Chem* **79**, 8185-8189 (2007). |
| 1281 | 127 | Liu, H.-B., Zhong, H., Karpowicz, N., Chen, Y. & Zhang, X.-C. Terahertz spectroscopy and |
| 1282 | | imaging for defense and security applications. *Proc IEEE* **95**, 1514-1527 (2007). |
| 1283 | 128 | Golightly, R. S., Doering, W. E. & Natan, M. J. Surface-enhanced Raman spectroscopy and |
| 1284 | | homeland security: a perfect match?  (ACS Nano, 2009). |
| 1285 | 129 | Sattlecker, M., Stone, N., Smith, J. & Bessant, C. Assessment of robustness and transferability |
| 1286 | | of classification models built for cancer diagnostics using Raman spectroscopy. *J Raman* |
| 1287 | | *Spectrosc* **42**, 897-903 (2011). |
| 1288 | 130 | Isabelle, M. *et al.* Multi-centre Raman spectral mapping of oesophageal cancer tissues: a |
| 1289 | | study to assess system transferability. *Faraday Discuss* **187**, 87-103 (2016). |
| 1290 | 131 | Guo, S. *et al.* Towards an improvement of model transferability for Raman spectroscopy in |
| 1291 | | biological applications. *Vib Spectrosc* **91**, 111-118 (2017). |
| 1292 | 132 | Luo, X. *et al.* Calibration transfer across near infrared spectrometers for measuring |
| 1293 | | hematocrit in the blood of grazing cattle. *Journal of Near Infrared Spectroscopy* **25**, 15-25 |
| 1294 | | (2017). |
| 1295 | 133 | Vaughan, A. A. *et al.* Liquid chromatography–mass spectrometry calibration transfer and |
| 1296 | | metabolomics data fusion. *Anal Chem* **84**, 9848-9857 (2012). |
| 1297 | 134 | Rodriguez, J. D., Westenberger, B. J., Buhse, L. F. & Kauffman, J. F. Standardization of Raman |
| 1298 | | spectra for transfer of spectral libraries across different instruments. *Analyst* **136**, 4232-4240 |
| 1299 | | (2011). |
| 1300 | 135 | de Andrade, E. W., de Lelis Medeiros de Morais, C., Lopes da Costa, F. S., de Lima, G. & |
| 1301 | | Michell, K. A Multivariate Control Chart Approach for Calibration Transfer between NIR |
| 1302 | | Spectrometers for Simultaneous Determination of Rifampicin and Isoniazid in |
| 1303 | | Pharmaceutical Formulation. *Curr Anal Chem* **14**, 488-494 (2018). |
| 1304 | 136 | Yu, B., Ji, H. & Kang, Y. Standardization of near infrared spectra based on multi-task learning. |
| 1305 | | *Spectroscopy Letters* **49**, 23-29 (2016). |
| 1306 | 137 | Ni, L., Han, M., Luan, S. & Zhang, L. Screening wavelengths with consistent and stable signals |
| 1307 | | to realize calibration model transfer of near infrared spectra. *Spectrochimica Acta Part A:* |
| 1308 | | *Molecular and Biomolecular Spectroscopy* (2018). |
| 1309 | 138 | Hu, R. & Xia, J. Calibration transfer of near infrared spectroscopy based on DS algorithm. in |
| 1310 | | *Electric Information and Control Engineering (ICEICE), 2011 International Conference on.* |
| 1311 | | 3062-3065 (IEEE). |
| 1312 | 139 | Forina, M. *et al.* Transfer of calibration function in near-infrared spectroscopy. *Chemom* |
| 1313 | | *Intellig Lab Syst* **27**, 189-203 (1995). |
| 1314 | 140 | Xiao, H. *et al.* Comparison of benchtop Fourier-transform (FT) and portable grating scanning |
| 1315 | | spectrometers for determination of total soluble solid contents in single grape berry (Vitis |
| 1316 | | vinifera L.) and calibration transfer. *Sensors* **17**, 2693 (2017). |
| 1317 | 141 | Yahaya, O., MatJafri, M., Aziz, A. & Omar, A. Visible spectroscopy calibration transfer model |
| 1318 | | in determining pH of Sala mangoes. *Journal of Instrumentation* **10**, T05002 (2015). |
| 1319 | 142 | Bin, J., Li, X., Fan, W., Zhou, J.-h. & Wang, C.-w. Calibration transfer of near-infrared |
| 1320 | | spectroscopy by canonical correlation analysis coupled with wavelet transform. *Analyst* **142**, |
| 1321 | | 2229-2238 (2017). |
| 1322 | 143 | Monakhova, Y. B. & Diehl, B. W. Transfer of multivariate regression models between high‐ |
| 1323 | | resolution NMR instruments: application to authenticity control of sunflower lecithin. |
| 1324 | | *Magnetic Resonance in Chemistry* **54**, 712-717 (2016). |

1325 144 Zuo, Q., Xiong, S., Chen, Z.-P., Chen, Y. & Yu, R.-Q. A novel calibration strategy based on
1326   background correction for quantitative circular dichroism spectroscopy. *Talanta* **174**, 320-
1327   324 (2017).
1328 145 Koehler IV, F. W., Small, G. W., Combs, R. J., Knapp, R. B. & Kroutil, R. T. Calibration transfer
1329   algorithm for automated qualitative analysis by passive Fourier transform infrared
1330   spectrometry. *Anal Chem* **72**, 1690-1698 (2000).
1331 146 Rodrigues, R. R. *et al.* Evaluation of calibration transfer methods using the ATR-FTIR
1332   technique to predict density of crude oil. *Chemom Intellig Lab Syst* **166**, 7-13 (2017).
1333 147 Wang, Y., Veltkamp, D. J. & Kowalski, B. R. Multivariate instrument standardization. *Anal*
1334   *Chem* **63**, 2750-2756 (1991).
1335 148 Brouckaert, D., Uyttersprot, J.-S., Broeckx, W. & De Beer, T. Calibration transfer of a Raman
1336   spectroscopic quantification method for the assessment of liquid detergent compositions
1337   from at-line laboratory to in-line industrial scale. *Talanta* **179**, 386-392 (2018).
1338 149 Andrade, E. V., Morais, C. d. L. M., Costa, F. S. L. & Lima, K. M. G. A Multivariate Control
1339   Chart Approach for Calibration Transfer between NIR Spectrometers for Simultaneous
1340   Determination of Rifampicin and Isoniazid in Pharmaceutical Formulation. *Curr Anal Chem*
1341   **14**, 1-7 (2018).
1342 150 Zamora-Rojas, E., Pérez-Marín, D., De Pedro-Sanz, E., Guerrero-Ginel, J. & Garrido-Varo, A.
1343   Handheld NIRS analysis for routine meat quality control: Database transfer from at-line
1344   instruments. *Chemom Intellig Lab Syst* **114**, 30-35 (2012).
1345 151 Panchuk, V., Kirsanov, D., Oleneva, E., Semenov, V. & Legin, A. Calibration transfer between
1346   different analytical methods. *Talanta* **170**, 457-463 (2017).
1347 152 de Morais, C. d. L. M. & de Lima, K. M. G. Determination and analytical validation of
1348   creatinine content in serum using image analysis by multivariate transfer calibration
1349   procedures. *Anal Methods* **7**, 6904-6910 (2015).
1350 153 Khaydukova, M. *et al.* Multivariate calibration transfer between two different types of
1351   multisensor systems. *Sensors Actuators B: Chem* **246**, 994-1000 (2017).
1352 154 Barreiro, P. *et al.* Calibration Transfer Between Portable and Laboratory NIR
1353   Spectrophotometers. *Acta Hortic* (2008).
1354 155 Sulub, Y., LoBrutto, R., Vivilecchia, R. & Wabuyele, B. W. Content uniformity determination
1355   of pharmaceutical tablets using five near-infrared reflectance spectrometers: a process
1356   analytical technology (PAT) approach using robust multivariate calibration transfer
1357   algorithms. *Anal Chim Acta* **611**, 143-150 (2008).
1358 156 Zhang, L., Small, G. W. & Arnold, M. A. Multivariate calibration standardization across
1359   instruments for the determination of glucose by Fourier transform near-infrared
1360   spectrometry. *Anal Chem* **75**, 5905-5915 (2003).
1361 157 Martens, H., Høy, M., Wise, B. M., Bro, R. & Brockhoff, P. B. Pre‐whitening of data by
1362   covariance‐weighted pre‐processing. *J Chemom* **17**, 153-165 (2003).
1363 158 Feudale, R. N. *et al.* Transfer of multivariate calibration models: a review. *Chemom Intellig*
1364   *Lab Syst* **64**, 181-192 (2002).
1365 159 Woody, N. A., Feudale, R. N., Myles, A. J. & Brown, S. D. Transfer of multivariate calibrations
1366   between four near-infrared spectrometers using orthogonal signal correction. *Anal Chem* **76**,
1367   2595-2600 (2004).
1368 160 Greensill, C., Wolfs, P., Spiegelman, C. & Walsh, K. Calibration transfer between PDA-based
1369   NIR spectrometers in the NIR assessment of melon soluble solids content. *Appl Spectrosc* **55**,
1370   647-653 (2001).
1371 161 Sjöblom, J., Svensson, O., Josefson, M., Kullberg, H. & Wold, S. An evaluation of orthogonal
1372   signal correction applied to calibration transfer of near infrared spectra. *Chemom Intellig Lab*
1373   *Syst* **44**, 229-244 (1998).
1374 162 Andrews, D. T. & Wentzell, P. D. Applications of maximum likelihood principal component
1375   analysis: incomplete data sets and calibration transfer. *Anal Chim Acta* **350**, 341-352 (1997).

1376 163 Bouveresse, E., Massart, D. & Dardenne, P. Calibration transfer across near-infrared
1377 spectrometric instruments using Shenk's algorithm: effects of different standardisation
1378 samples. *Anal Chim Acta* **297**, 405-416 (1994).
1379 164 Shenk, J. S. & Westerhaus, M. O. Populations structuring of near infrared spectra and
1380 modified partial least squares regression. *Crop Sci.* **31**, 1548-1555 (1991).
1381 165 Paatero, P. & Tapper, U. Positive matrix factorization: A non‐negative factor model with
1382 optimal utilization of error estimates of data values. *Environmetrics* **5**, 111-126 (1994).
1383 166 Xie, Y. & Hopke, P. K. Calibration transfer as a data reconstruction problem. *Anal Chim Acta*
1384 **384**, 193-205 (1999).
1385 167 Goodacre, R. *et al.* On mass spectrometer instrument standardization and interlaboratory
1386 calibration transfer using neural networks. *Anal Chim Acta* **348**, 511-532 (1997).
1387 168 Chen, W.-R., Bin, J., Lu, H.-M., Zhang, Z.-M. & Liang, Y.-Z. Calibration transfer via an extreme
1388 learning machine auto-encoder. *Analyst* **141**, 1973-1980 (2016).
1389 169 Hu, Y., Peng, S., Bi, Y. & Tang, L. Calibration transfer based on maximum margin criterion for
1390 qualitative analysis using Fourier transform infrared spectroscopy. *Analyst* **137**, 5913-5918
1391 (2012).
1392 170 Fan, W., Liang, Y., Yuan, D. & Wang, J. Calibration model transfer for near-infrared spectra
1393 based on canonical correlation analysis. *Anal Chim Acta* **623**, 22-29 (2008).
1394 171 Wang, Z., Dean, T. & Kowalski, B. R. Additive background correction in multivariate
1395 instrument standardization. *Anal Chem* **67**, 2379-2385 (1995).
1396 172 Kennard, R. W. & Stone, L. A. Computer Aided Design of Experiments. *Technometrics* **11**,
1397 137-148 (1969).
1398 173 Palonpon, A. F. *et al.* Raman and SERS microscopy for molecular imaging of live cells. *Nat*
1399 *Protoc* **8**, 677 (2013).
1400 174 Witze, E. S., Old, W. M., Resing, K. A. & Ahn, N. G. Mapping protein post-translational
1401 modifications with mass spectrometry. *Nat Methods* **4**, 798 (2007).
1402 175 Aebersold, R. & Mann, M. Mass spectrometry-based proteomics. *Nature* **422**, 198 (2003).
1403 176 Pence, I. & Mahadevan-Jansen, A. Clinical instrumentation and applications of Raman
1404 spectroscopy. *Chem Soc Rev* **45**, 1958-1979 (2016).
1405 177 Ibrahim, O. *et al.* Improved protocols for pre-processing Raman spectra of formalin fixed
1406 paraffin preserved tissue sections. *Anal Methods* **9**, 4709-4717 (2017).
1407 178 Tfayli, A. *et al.* Digital dewaxing of Raman signals: discrimination between nevi and
1408 melanoma spectra obtained from paraffin-embedded skin biopsies. *Appl Spectrosc* **63**, 564-
1409 570 (2009).
1410 179 Byrne, H. J., Knief, P., Keating, M. E. & Bonnier, F. Spectral pre and post processing for
1411 infrared and Raman spectroscopy of biological tissues and cells. *Chem Soc Rev* **45**, 1865-1878
1412 (2016).
1413 180 Meade, A. D. *et al.* Studies of chemical fixation effects in human cell lines using Raman
1414 microspectroscopy. *Anal Bioanal Chem* **396**, 1781-1791 (2010).
1415 181 Baker, M. J. *et al.* Developing and understanding biofluid vibrational spectroscopy: a critical
1416 review. *Chem Soc Rev* **45**, 1803-1818 (2016).
1417 182 Bonifacio, A., Cervo, S. & Sergo, V. Label-free surface-enhanced Raman spectroscopy of
1418 biofluids: fundamental aspects and diagnostic applications. *Anal Bioanal Chem* **407**, 8265-
1419 8277 (2015).
1420 183 Mitchell, A. L., Gajjar, K. B., Theophilou, G., Martin, F. L. & Martin-Hirsch, P. L. Vibrational
1421 spectroscopy of biofluids for disease screening or diagnosis: translation from the laboratory
1422 to a clinical setting. *J Biophotonics* **7**, 153-165 (2014).
1423 184 Lovergne, L. *et al.* Biofluid infrared spectro-diagnostics: pre-analytical considerations for
1424 clinical applications. *Faraday Discuss* **187**, 521-537 (2016).
1425 185 Bonifacio, A. *et al.* Surface-enhanced Raman spectroscopy of blood plasma and serum using
1426 Ag and Au nanoparticles: a systematic study. *Anal Bioanal Chem* **406**, 2355-2365 (2014).

1427 186 Paraskevaidi, M., Martin-Hirsch, P. L. & Martin, F. L. ATR-FTIR Spectroscopy Tools for Medical
1428      Diagnosis and Disease Investigation. *Springer* (2018).
1429 187 Mitchell, B. L., Yasui, Y., Li, C. I., Fitzpatrick, A. L. & Lampe, P. D. Impact of freeze-thaw cycles
1430      and storage time on plasma samples used in mass spectrometry based biomarker discovery
1431      projects. *Cancer Inform* **1** (2005).
1432 188 Glassford, S. E., Byrne, B. & Kazarian, S. G. Recent applications of ATR FTIR spectroscopy and
1433      imaging to proteins. *Biochim Biophys Acta* **1834**, 2849-2858 (2013).
1434 189 Kundu, J., Le, F., Nordlander, P. & Halas, N. J. Surface enhanced infrared absorption (SEIRA)
1435      spectroscopy on nanoshell aggregate substrates. *Chem Phys Lett* **452**, 115-119 (2008).
1436 190 Jones, S., Carley, S. & Harrison, M. An introduction to power and sample size estimation.
1437      *Emergency Medicine Journal* **20**, 453-458 (2003).
1438 191 Beebe, K. R., Pell, R. J. & Seasholtz, M. B. *Chemometrics: a practical guide*. Vol. **4** (Wiley New
1439      York, 1998).
1440 192 Pavia, D. L., Lampman, G. M., Kriz, G. S. & Vyvyan, J. A. *Introduction to spectroscopy*.
1441      (Cengage Learning, 2008).
1442 193 Savitzky, A. & Golay, M. J. Smoothing and differentiation of data by simplified least squares
1443      procedures. *Anal Chem* **36**, 1627-1639 (1964).
1444 194 Geladi, P., MacDougall, D. & Martens, H. Linearization and scatter-correction for near-
1445      infrared reflectance spectra of meat. *Appl Spectrosc* **39**, 491-500 (1985).
1446 195 Barnes, R., Dhanoa, M. S. & Lister, S. J. Standard normal variate transformation and de-
1447      trending of near-infrared diffuse reflectance spectra. *Appl Spectrosc* **43**, 772-777 (1989).
1448 196 Brereton, R. G. *Chemometrics: data analysis for the laboratory and chemical plant*. (John
1449      Wiley & Sons, 2003).
1450 197 Hastie, T., Tibshirani, R. & Friedman, J. The elements of statistical learning 2nd edition (New
1451      York: Springer, 2009).
1452 198 Bro, R. & Smilde, A. K. Principal component analysis. *Anal Methods* **6**, 2812-2831 (2014).
1453 199 Martin, F. L. *et al.* Identifying variables responsible for clustering in discriminant analysis of
1454      data from infrared microspectroscopy of a biological sample. *J Comput Biol* **14**, 1176-1184
1455      (2007).
1456 200 Martens, H. & Martens, M. Modified Jack-knife estimation of parameter uncertainty in
1457      bilinear modelling by partial least squares regression (PLSR). *Food quality and preference* **11**,
1458      5-16 (2000).
1459 201 Rousseeuw, P. J. & Hubert, M. Robust statistics for outlier detection. *Wiley Interdisciplinary
1460      Reviews: Data Mining and Knowledge Discovery* **1**, 73-79 (2011).
1461 202 Jiang, F., Liu, G., Du, J. & Sui, Y. Initialization of K-modes clustering using outlier detection
1462      techniques. *Inf Sci* **332**, 167-183 (2016).
1463 203 Domingues, R., Filippone, M., Michiardi, P. & Zouaoui, J. A comparative evaluation of outlier
1464      detection algorithms: Experiments and analyses. *Pattern Recognit* **74**, 406-421 (2018).
1465 204 Bakeev, K. A. *Process analytical technology: spectroscopic tools and implementation
1466      strategies for the chemical and pharmaceutical industries*. (John Wiley & Sons, 2010).
1467 205 Kuligowski, J., Quintás, G., Herwig, C. & Lendl, B. A rapid method for the differentiation of
1468      yeast cells grown under carbon and nitrogen-limited conditions by means of partial least
1469      squares discriminant analysis employing infrared micro-spectroscopic data of entire yeast
1470      cells. *Talanta* **99**, 566-573 (2012).
1471 206 Morais, C. L. & Lima, K. M. Comparing unfolded and two-dimensional discriminant analysis
1472      and support vector machines for classification of EEM data. *Chemom Intellig Lab Syst* (2017).
1473 207 Dixon, S. J. & Brereton, R. G. Comparison of performance of five common classifiers
1474      represented as boundary methods: Euclidean Distance to Centroids, Linear Discriminant
1475      Analysis, Quadratic Discriminant Analysis, Learning Vector Quantization and Support Vector
1476      Machines, as dependent on data structure. *Chemom Intellig Lab Syst* **95**, 1-17 (2009).

1477 208 Brereton, R. G. & Lloyd, G. R. Partial least squares discriminant analysis: taking the magic
1478    away. *J Chemom* **28**, 213-225 (2014).
1479 209 Cover, T. & Hart, P. Nearest neighbor pattern classification. *IEEE Trans Inf Theory* **13**, 21-27
1480    (1967).
1481 210 Cortes, C. & Vapnik, V. Support-vector networks. *Mach Learn* **20**, 273-297 (1995).
1482 211 Abraham, A. Artificial neural networks. *handbook of measuring system design* (2005).
1483 212 Fawagreh, K., Gaber, M. M. & Elyan, E. Random forests: from early developments to recent
1484    advancements. *Systems Science & Control Engineering: An Open Access Journal* **2**, 602-609
1485    (2014).
1486 213 LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *nature* **521**, 436 (2015).
1487 214 Seasholtz, M. B. & Kowalski, B. The parsimony principle applied to multivariate calibration.
1488    *Anal Chim Acta* **277**, 165-177 (1993).
1489 215 Morais, C. L. & Lima, K. M. Principal Component Analysis with Linear and Quadratic
1490    Discriminant Analysis for Identification of Cancer Samples Based on Mass Spectrometry. *J*
1491    *Braz Chem Soc*, 31 (2017).
1492 216 Hibbert, D. B. Vocabulary of concepts and terms in chemometrics (IUPAC Recommendations
1493    2016). *Pure and Applied Chemistry* **88**, 407-443 (2016).
1494 217 McCall, J. Genetic algorithms for modelling and optimisation. *J Comput Appl Math* **184**, 205-
1495    222 (2005).
1496 218 Soares, S. F. C., Gomes, A. A., Araujo, M. C. U., Galvão Filho, A. R. & Galvão, R. K. H. The
1497    successive projections algorithm. *Trends Anal Chem* **42**, 84-98 (2013).
1498 219 Kamandar, M. & Ghassemian, H. Maximum relevance, minimum redundancy feature
1499    extraction for hyperspectral images. in *Electrical Engineering (ICEE), 2010 18th Iranian*
1500    *Conference on.* 254-259 (IEEE).

1501

# Supplementary Material 1


Additional results from pilot study
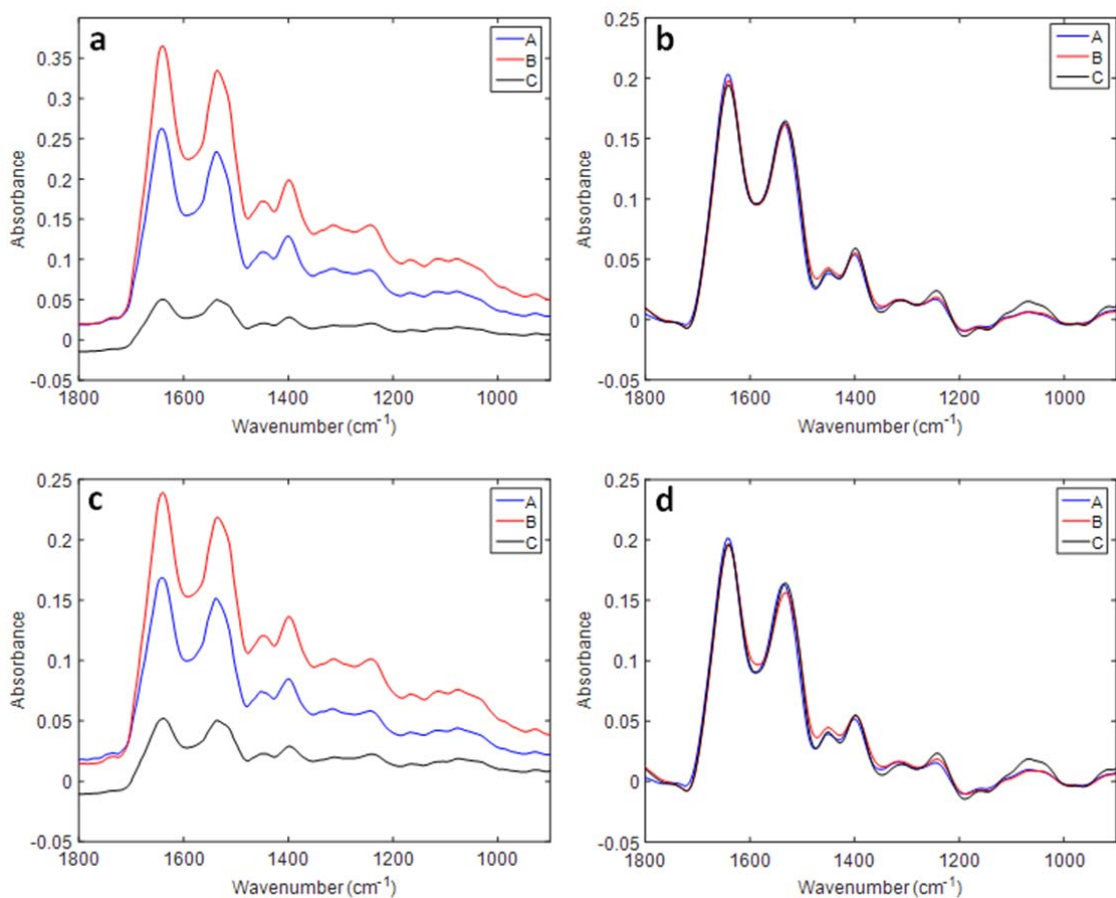
# A. Effect of different instruments



**Figure S1.** Average (a) raw and (b) pre-processed spectra for healthy controls samples; average (c) raw and (d) pre-processed spectra for cancer samples across three different instruments (A, B and C).
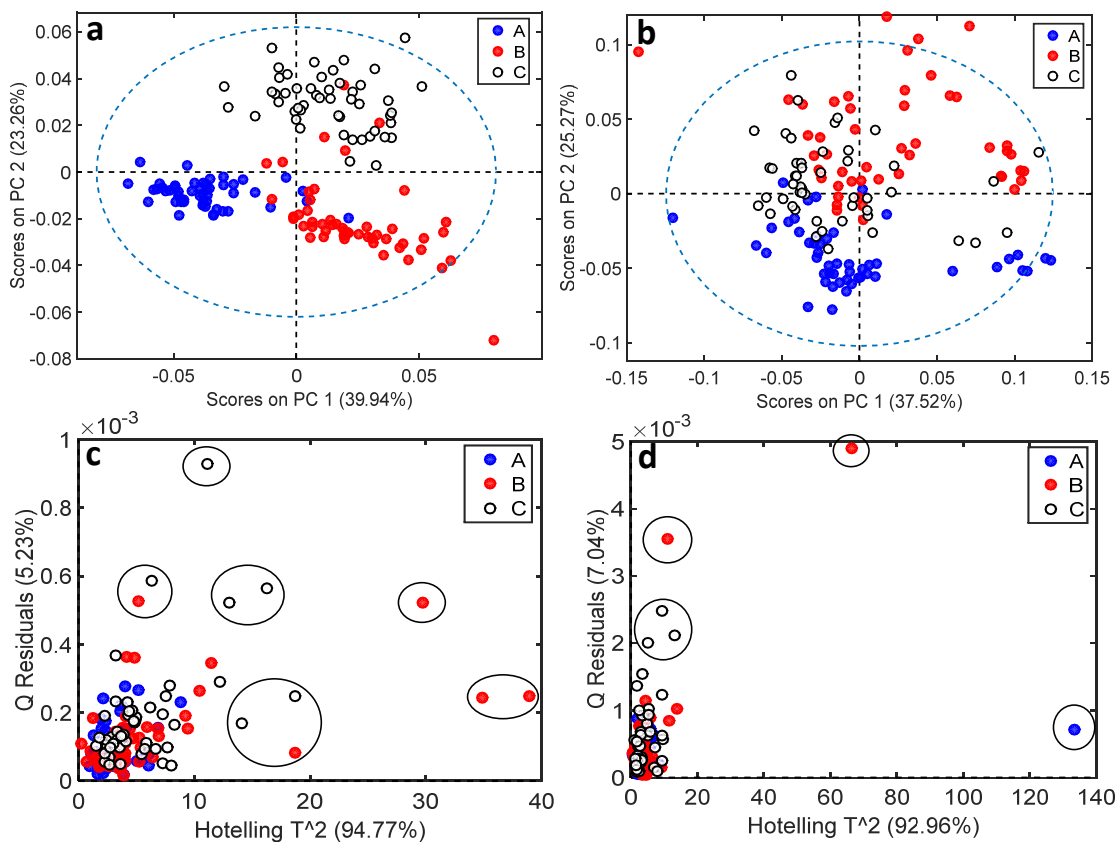
## A. Effect of different instruments



**Figure S2.** (a) PCA scores for healthy control samples according to the instrument used for spectra acquisition (A, B and C); (b) PCA scores for cancer samples according to the instrument used for spectra acquisition (A, B and C); (c) Hotelling $T^2$ *versus* Q residual test for healthy control samples according to the instrument used for spectra acquisition (A, B and C) based on a PCA using 5 PCs (94.77% cumulative variance); (d) Hotelling $T^2$ *versus* Q residual test for cancer samples according to the instrument used for spectra acquisition (A, B and C) based on a PCA using 5 PCs (92.96% cumulative variance). Circled samples in (c) and (d) indicate outliers removed. Confidence ellipse was 95%, depicted in blue in (a) and (b).

**Figure S3.** (a) PCA loadings for healthy control samples measured in different instruments (A, B and C); (b) PCA loadings for cancer samples measured in different instruments (A, B and C).

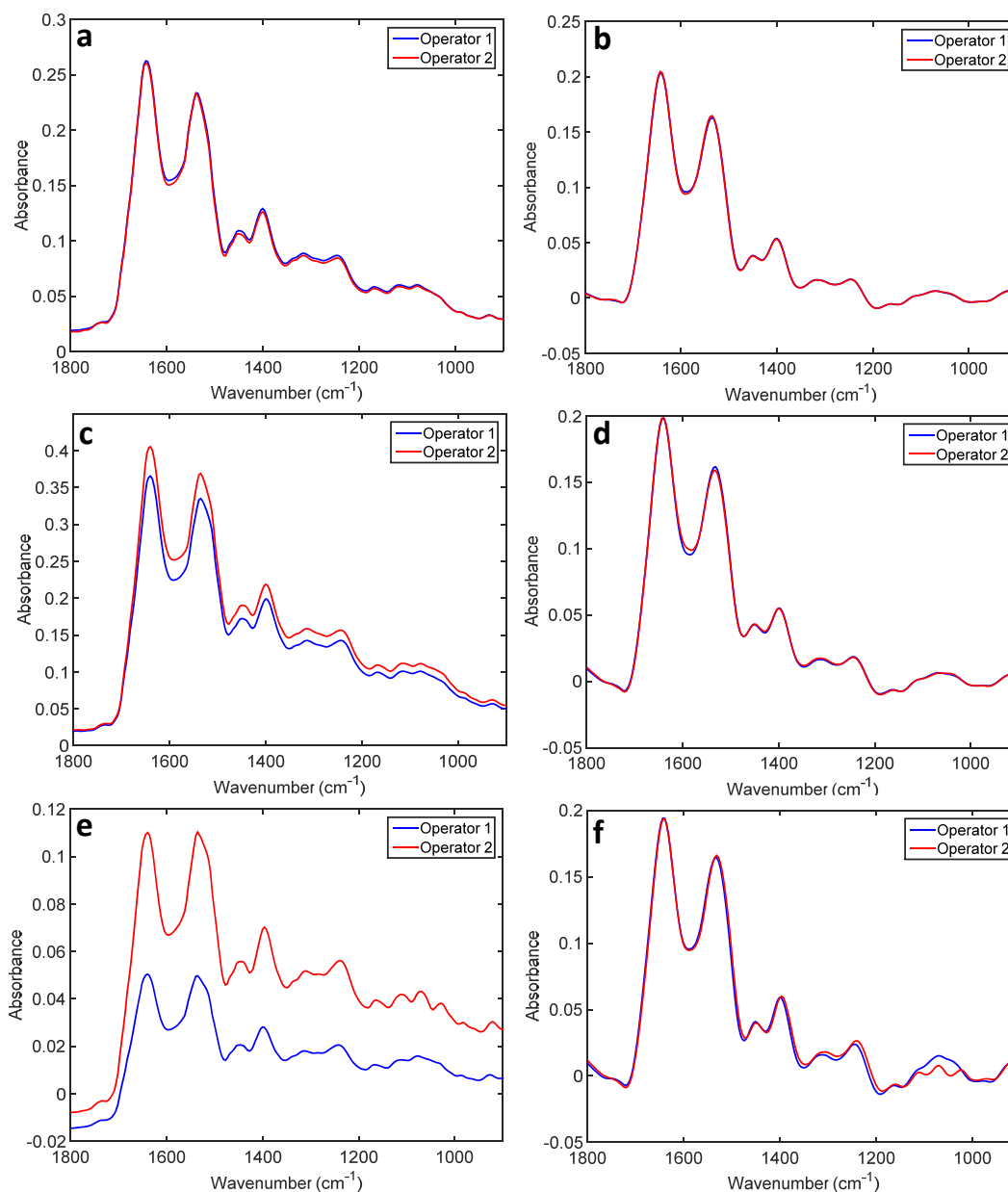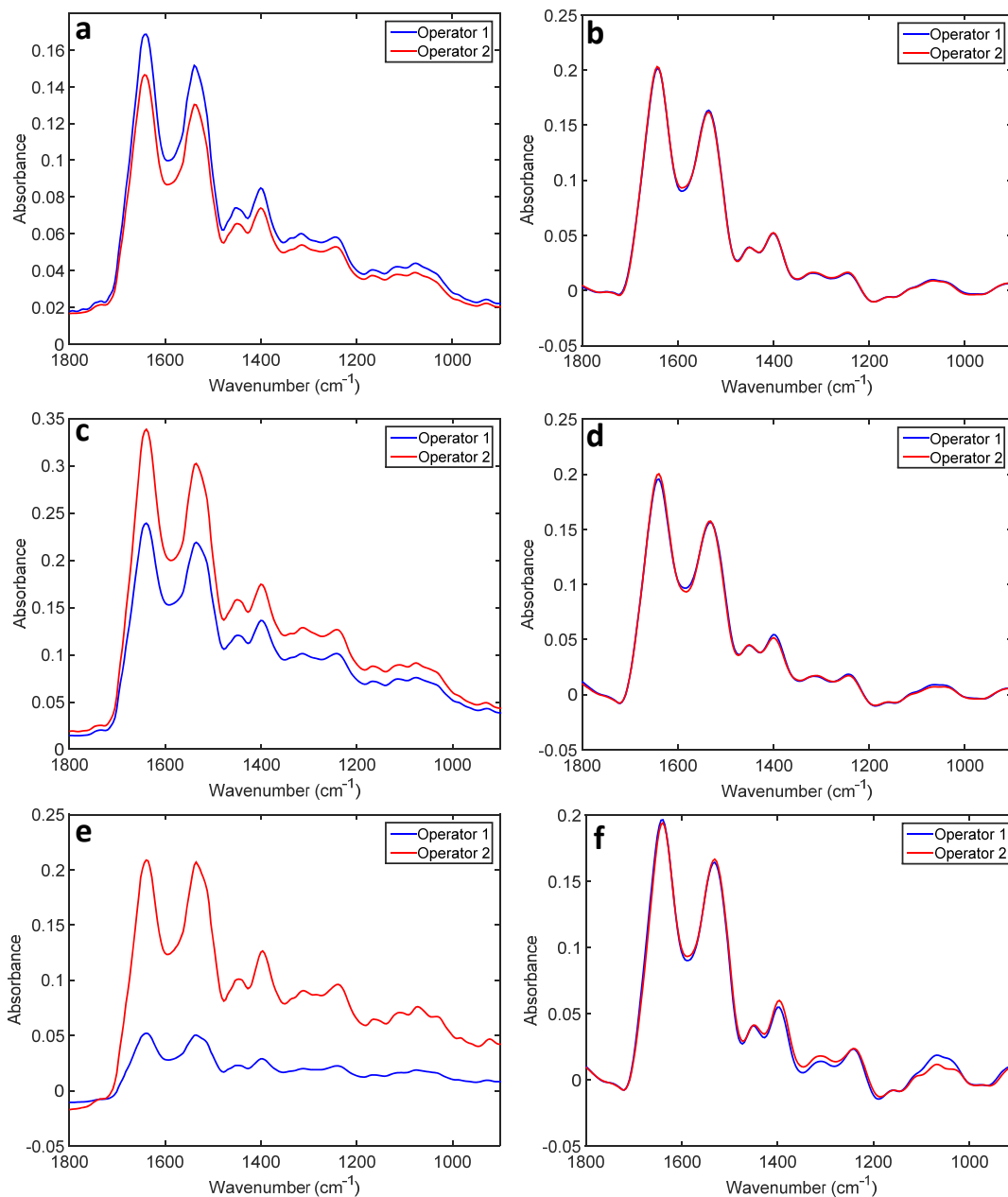# B. Effect of different operators



**Figure S4.** Average (a) raw and (b) pre-processed spectra for healthy control samples acquired with instrument A depending on the operator; average (c) raw and (d) pre-processed spectra for healthy control samples acquired with instrument B depending on the operator; average (e) raw and (f) pre-processed spectra for healthy control samples acquired with instrument C varying the operator.

# B. Effect of different operators



**Figure S5.** Average (a) raw and (b) pre-processed spectra for cancer samples acquired with instrument A depending on the operator; average (c) raw and (d) pre-processed spectra for cancer samples acquired with instrument B depending on the operator; average (e) raw and (f) pre-processed spectra for cancer samples acquired with instrument C depending on the operator.
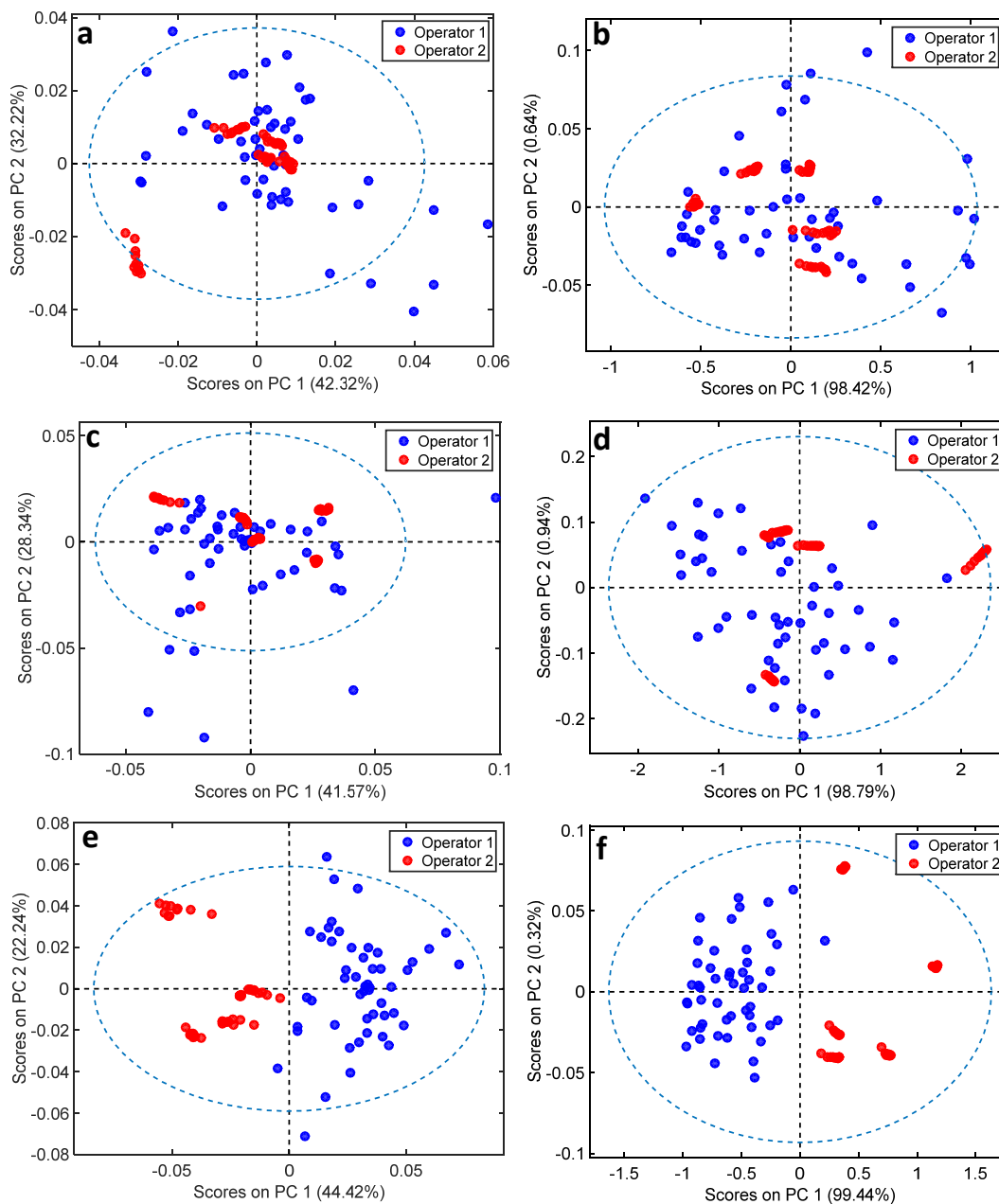
# B. Effect of different operators



**Figure S6.** PCA scores for (a) healthy control and (b) cancer samples acquired with instrument A depending on the operator; PCA scores for (c) healthy control and (d) cancer samples acquired with instrument B depending on the operator; PCA scores for (e) healthy control and (f) cancer samples acquired with instrument C depending on the operator. Confidence ellipse was 95%, depicted in blue

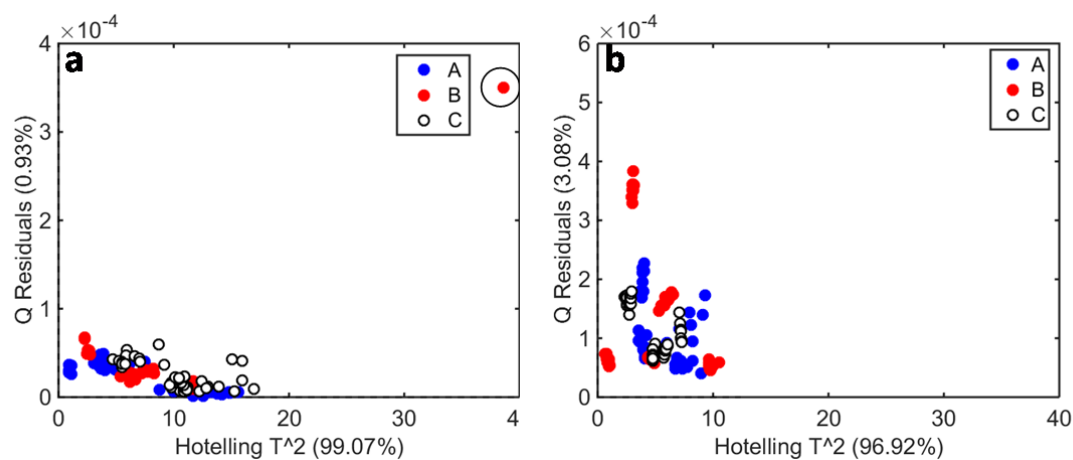## C. Effect of different instruments and operators



**Figure S7.** (a) Hotelling $T^2$ *versus* Q residual test based on a PCA using 8 PCs (99.07% cumulative variance) for healthy control samples depending on the instrument for spectra acquisition (A, B and C) used by Operator 2; (b) Hotelling $T^2$ *versus* Q residual test based on a PCA using 5 PCs (96.92% cumulative variance) for cancer samples depending on the instrument for spectra acquisition (A, B and C) used by Operator 2. Circled sample in a) indicates an outlier removed. The Hotelling $T^2$ *versus* Q residual test for Operator 1 is depicted in Fig. S2c-d.
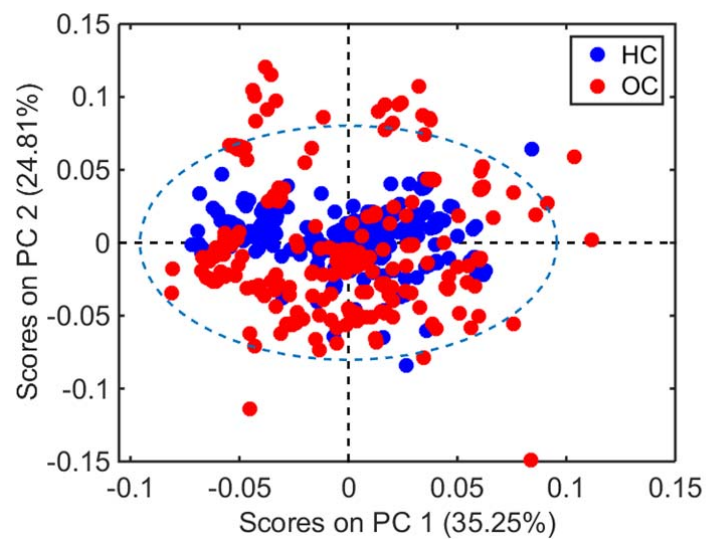
## D. Effect of different classes



**Figure S8.** PCA scores for healthy controls (HC) and ovarian cancer (OC) samples based on the spectra acquired by both operators (1 and 2) and by all instruments (A, B and C). Confidence ellipse at a 95% confidence level is depicted in blue

# Supplementary Method 1

Protocol for outliers detection

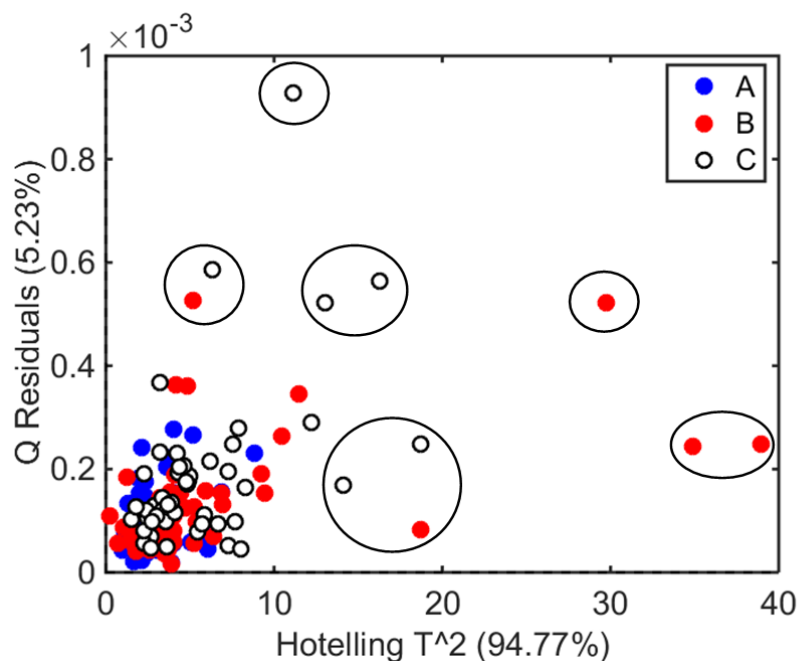# A. Outlier detection using Hotelling T$^2$ *versus* Q residuals test

**1$^{st}$ step:** Build a PCA model.

**2$^{nd}$ step:** Calculate Hotelling T$^2$ and Q residuals.

**3$^{rd}$ step:** Plot Hotelling T$^2$ *versus* Q residuals

**4$^{th}$ step:** Select the samples which are most distant to the plot origin (0,0) and remove them one at a time from the data set. This procedure can be performed manually after visual inspection or automatically by algorithms.

**Figure S1.** Hotelling T$^2$ *versus* Q residuals for healthy control samples (blood plasma) varying the instrument for spectra acquisition (A, B and C). PCA performed with 5 PCs (94.77% cumulative variance). Circled samples indicate outliers removed.

# B. Automatic outlier detection using MATLAB®

Algorithm link to download:

https://doi.org/10.6084/m9.figshare.7066613.v1

**1st step:** Add the .m files within the file downloaded to the path.

**2nd step:** Load the spectral data into MATLAB and organize all the spectra into a single matrix "X" containing each spectrum as a row.

**3rd step:** Perform an initial PCA model to determine the number of principal components (PCs) to work with.

**4th step:** Run the algorithm as follows:

```
Command Window
fx >> Xc = outlier(X,Npcs);
```

where "Xc" is the spectral matrix without outliers, "X" is the input spectral data, and "Npcs" the number of PCs for PCA.

**5th step:** Input optimization parameters:

```
Command Window
   >> Xc = outlier(X,Npcs);
   ------------------
   Select the Hotelling T2 threshold: 25
   Select the Q residuals threshold: 0.8e-03
fx Select the number of repetitions: 10
```

In this case, the algorithm will perform a PCA model 10 times removing one sample at a time that follows one of these criteria: Hotteling $T^2 > 25$ or Q residuals $> 0.8 \times 10^{-3}$. Then, these samples are automatic excluded from the new dataset (Xc). The list of excluded samples is also displayed in MATLAB. Example:

```
Command Window
    >> Xc = outlier(X,Npcs);
    ------------------
    Select the Hotelling T2 threshold: 25
    Select the Q residuals threshold: 0.8e-03
    Select the number of repetitions: 10
    ------------------
    Removed samples:

        97

        97

        77

       141


    ------------------
fx >> |
```