

1 **The metabolomic detection of lung cancer biomarkers in sputum**

2

3 Simon J. S. Cameron¹, Keir E. Lewis^{2,3}, Manfred Beckmann¹, Gordon G. Allison¹, Robin Ghosal², Paul D. Lewis³,
4 and Luis A. J. Mur^{1*}.

5

6 ¹Institute of Biological, Environmental and Rural Sciences, Edward Llywd Building, Penglais Campus,
7 Aberystwyth, Ceredigion, SY23 3FG, UK.

8 ²Department of Respiratory Medicine, Prince Phillip Hospital, Llanelli, SA14 8LY, UK.

9 ³College of Medicine, Swansea University, Swansea, SA2 8PP, UK.

10

11 **Running Title:** Sputum Metabolomics for Lung Cancer Biomarkers

12

13 *** Corresponding Authors:** Professor Luis A. J. Mur. Institute of Biological, Environmental and Rural Sciences,
14 Edward Llywd Building, Aberystwyth University, Penglais Campus, Aberystwyth, Ceredigion, SY23 3FG, Wales,
15 UK. Phone: 01970 622981. Fax: 01970 622350. Email Address: lum@aber.ac.uk

16

17

18

19

20

21

22

23

24

25

26

27

28

29 **Abstract**

30

31 **Objectives:** Developing screening and diagnosis methodologies based on novel biomarkers should allow for the
32 detection of the lung cancer (LC) and possibly at an earlier stage and thereby increase the effectiveness of clinical
33 interventions. Here, our primary objective was to evaluate the potential of spontaneous sputum as a source of
34 non-invasive metabolomic biomarkers for LC status.

35

36 **Materials and Methods:** Spontaneous sputum was collected and processed from 34 patients with suspected LC,
37 alongside 33 healthy controls. Of the 34 patients, 23 were subsequently diagnosed with LC (LC⁺, 16 NSCLC, six
38 SCLC, and one radiological diagnosis), at various stages of disease progression. The 67 samples were analysed
39 using flow infusion electrospray ion mass spectrometry (FIE-MS) and gas-chromatography mass spectrometry
40 (GC-MS).

41

42 **Results:** Principal component analysis identified negative mode FIE-MS as having the main separating power
43 between samples from healthy and LC. Discriminatory metabolites were identified using ANOVA and Random
44 Forest. Indications of potential diagnostic accuracy involved the use of receiver operating characteristic / area
45 under the curve (ROC/AUC) analyses. This approach identified metabolites changes that were only observed
46 with LC. Metabolites with AUC values of greater than 0.8 which distinguished between LC⁺/LC⁻ binary
47 classifications were identified and included Ganglioside GM1 which has previously been linked to LC.

48

49 **Conclusion:** This study indicates that metabolomics based on sputum can yield metabolites that can be used as
50 a diagnostic and/or discriminator tool. These could aid clinical intervention and targeted diagnosis of LC within
51 an 'at risk' LC⁻ population group. The use of sputum as a non-invasive source of metabolite biomarkers may aid
52 in the development of an at-risk population screening programme for lung cancer or enhanced clinical diagnostic
53 pathways.

54

55 **Key Words:** Lung cancer, metabolomics, biomarkers, sputum, polyamines, gangliosides

56

57 **1 Introduction**

58

59 Lung cancer (LC) is the most prevalent cancer in the world; responsible for 1.3 million deaths annually [1]. The
60 last 30 years has seen little improvement in the overall five year survival rate for LC; with only 15% of patients
61 living for at least five years after their initial diagnosis [2]. These relatively poor survival rates are primarily a
62 result of the late detection of a malignancy; reducing the success of clinical interventions. Clinicians currently
63 rely on three main tools for LC diagnosis: X-ray, computerised tomography (CT) scans, and bronchoscopy. These
64 methods have improved our ability to detect lung cancer, but have nevertheless failed to improve the rate of
65 early LC detection [3]. Another aspect of this poor early detection is the association of LC with smoking, which
66 masks some of the disease's early symptoms, which has been linked to approximately 90% of LC tumours [4].

67

68 An alternative screening methodology to radiography, which is currently the most widely used approach, is the
69 utilisation of molecular markers, both genetic and metabolomic, in biofluids. For example, microRNAs have been
70 suggested as biomarkers for NSCLC in sputum [5], plasma [6], and serum [7]. Previous work by members of this
71 research group has demonstrated that chemometric analysis combined with Fourier transform infrared
72 spectroscopy is a non-invasive approach that allows for the discrimination of LC positive patients. This
73 demonstrated that sputum could be used as a non-invasive source of biomarkers for LC [8]. However, analysis
74 of mid-IR spectra only provides information on broad changes in classes of chemicals, and has a poor ability to
75 resolve changes to particular chemicals. By comparison, metabolite profiling based on sample screening using
76 Mass Spectrometry (MS) can resolve changes in individual chemicals and thus, could more readily identify
77 biomarkers linked to LC detection.

78

79 The aim of this study was to employ MS metabolomic profiling to identify clinically relevant biomarkers in
80 sputum that could be used for detect LC (diagnosis) as well as provide some pathophysiological insights based
81 on the characteristics of the chemical biomarkers. We utilised two MS approaches in this study, Gas
82 Chromatography MS (GC-MS) and Flow Infusion Electrospray MS (FIE-MS). Our rationale for this approach is that
83 both MS technologies are widely used in biomarker discovery, but have differing levels of sensitivities and
84 different approaches in regards to sample preparation and analysis. For example, GC-MS requires chemical

85 derivatization of sample metabolites prior to analysis whilst FIE-MS requires no pre-treatment [9]. Although, our
86 study employed both univariate and multivariate approaches our study sought to conform to the demands of
87 the TRIPOD (The Transparent Reporting of a multivariable prediction model for Individual Prognosis Or
88 Diagnosis) Statement by adhering to the recommended checklist [10]. We employed assessments of diagnostic
89 accuracy based on receiver operating characteristic (ROC)/ Area under the Curve (AUC) that suggest that our
90 approach could be used in clinical context to inform the detection of the disease. To the best of our knowledge,
91 metabolomic profiles have not been reported using sputum as a biofluid from clinical patients. Thus, beyond,
92 the detection of biomarkers, a description of the LC sputum metabolome offers a novel insight into the
93 pathology of LC.

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114 **2 Materials and Methods**

115

116 **2.1 Ethics Statement**

117 The MedLung observational study (UKCRN ID 4682) received loco-regional ethical approval from the Hywel Dda
118 Health Board (05/WMW01/75). Written informed consent was obtained from all participants at least 24 hours
119 before sampling, at a previous clinical appointment, and all data was link anonymised before analysis.

120

121 **2.2 Study Design**

122 This study aimed to compare the metabolomes of three groups of sputum samples. Two sets of sputum samples
123 were obtained from patients referred to the access LC clinic at the Prince Phillip Hospital, Wales, UK; a site of
124 primary care. Lung cancer status was subsequently assessed as part of the Medlung observational study (UKCRN
125 ID 4682) and patients were classified as either LC⁺ or diagnosed with another pulmonary disease (LC⁻) based on
126 histological assessments of sputum bronchoscopy derived samples (Table S1). Metadata including
127 comorbidities, smoking history and drug history are given in Table S1. Additionally, spontaneous sputum
128 samples were collected from staff members at Swansea University who had no previous history of cancer or
129 lung disease, other than asthma. These non-clinical samples were designated as a control (CON) group. The
130 design extensively exploited pair-wise analyses between LC⁺ and LC⁻ groups and the CON group. As this project
131 was seen as a pilot project, no external validation set, comprising, for example, testing on another set of patients
132 samples was used. Further, the danger of over-fitting the derived data was reduced through the extensive use
133 of simple two-way ANOVA in our pairwise comparisons. Sampling occurred between 2012-2013 to align with
134 the Medlung study timeline and this, rather than an *a priori* design target, governed the number of samples
135 analysed.

136

137 **2.3 Patient Recruitment and Sampling**

138 Spontaneous sputum was collected from referrals to our rapid access LC clinic at the Prince Phillip Hospital,
139 Wales, UK or volunteers from staff members at Swansea University. No *a priori* criteria were applied to the
140 selection of patients or volunteers other than their ability to produce sputum. Patients were asked to cough
141 into sterile, 50 mL polypropylene tubes (Greiner Bio-One Ltd, UK) prior to bronchoscopy, to a total volume of

142 2-3 mL. A 100 μ L aliquot of all samples, including the CON group, was taken to create a second pellet that was
143 subsequently formalin fixed and wax embedded prior to sectioning and staining with haemotoxylin and eosin
144 (H&E). To confirm samples were of bronchial origin, H&E stained sections were assessed by a consultant
145 histopathologist for presence of bronchial epithelial cells. Histological assessments of the LC⁺ class allowed the
146 recording of LC type and stage. Thus, NSCLC classifications were obtained for sixteen samples and six were SCLC.
147 Only in one case (LC06) was no classification obtained. Within the NSCLC samples, seven could be sub-classified
148 as adenocarcinoma type and five squamous cell types. Considering the LC⁻ classified samples, three were
149 diagnosed with chronic obstructive pulmonary disease (COPD) and two with pneumonia. Amongst the LC⁺ group,
150 only two (LC07, LC20) were diagnosed with COPD which could be considered a LC co-morbidity and none with
151 pneumonia.

152

153 **2.4 Processing of Raw Sputum**

154 In line Raw sputum samples were frozen at -80 °C and defrosted in ice for approximately two hours when
155 required. Sputum cells were isolated by adding 0.5 mL of a working solution of dithiothreitol (DTT), made up by
156 adding 2.5 g of DTT to 31 mL of 30% aqueous methanol, and 5 mL of 30% aqueous methanol. The samples were
157 then placed on a vortex mixer for 15 minutes and underwent centrifugation at 1,800 x g for 10 minutes. The
158 supernatant was removed and the pellet used in subsequent metabolomic profiling.

159

160 **2.5 Flow Infusion Electrospray Mass Spectrometry (FIE-MS)**

161 After processing, 20 μ L of the sputum pellet was added to 20 μ L of ultrapure water and 40 μ L of ice-cold HPLC
162 grade acetone. Samples were vortex mixed for five seconds, cooled on ice for 30 minutes, and then underwent
163 centrifugation at 11,000 x g for five minutes. After centrifugation, 50 μ L of the supernatant was removed and
164 250 μ L of 70% methanol (made up using HPLC grade methanol and ultrapure water) was added. Glass vials were
165 capped and analysed in random order on a LTQ linear ion trap (Thermo Electron Corporation). Data were
166 acquired in alternating positive and negative ionization modes over 4 scan ranges (15–110, 100–220, 210–510,
167 and 500–1200 m/z), with an acquisition time of five minutes. The resulting mass spectrum was the mean of 20
168 scans about the apex of the infusion profile.

169

170 **2.6 Gas Chromatography Mass Spectrometry (GC-MS)**

171 The sputum pellet was processed as described in section 2.4 and 50 μL of the supernatant after centrifugation
172 removed and dried using a DNA SpeedVac (Savant, USA) at 40°C. After removal of all liquid, 30 μL of a 20 mg/ml
173 solution of methoxyamine dissolved in pyridine was added and each sample was transferred to a 11 mm
174 diameter glass GC vials which were capped with Teflon crimp caps and incubated at 90°C for 15 minutes. After
175 cap removal, 20 μL of N,O-Bis(trimethylsilyl)trifluoroacetamide (BSTFA) was added to the sample, alongside 5
176 μL of an alkane standard mix. This mixture comprised of C₁₀, C₁₃, C₁₅, C₁₈, C₁₉, C₂₃, C₂₈, C₃₂ and C₃₆ alkanes dissolved
177 in pyridine each at a concentration of 2 $\mu\text{L}/\text{mL}$ (for alkanes liquid at room temperature) or 2 mg/mL (for alkanes
178 solid at room temperature). The vials were recapped and incubated at 90°C for 15 minutes. Samples were
179 analysed by duplicate injection on a 6890N GC linked to a 5973N mass analyser and a 7683 auto-sampler (Agilent
180 Technologies) fitted with a Thermo Scientific TR-Pesticide II fused silica column (30m x 0.25 mm ID x 0.25 μm
181 film thickness). Helium carrier gas was supplied at a constant flow rate of 1 mL per minute and following the
182 injection of 1 μL of sample the GC oven was held at 80°C for three minutes, increased to 280°C at a rate of 15°C
183 per minute, and then to 330°C at a rate of 50°C per minute. The inlet temperature was 280°C and samples were
184 split with a 2:1 split ratio. The temperature of the MS transfer line was 330°C.

185

186 **2.7 Accurate Mass Determination**

187 Selected discriminatory nominal mass signals were investigated further by targeted nano-flow Fourier
188 Transform-Ion Cyclotron Resonance Ultra-Mass-Spectrometry (FT-ICR-MS) using TriVersa NanoMate (Advion
189 BioSciences Ltd) on a LTQ-FT-ULTRA (Thermo Scientific) to obtain ultra-high accurate mass information and MSn
190 ion-trees [11]. Resulting accurate mass values were used to interrogate the Human Metabolome Database [12].
191 Based on an accuracy of 1 ppm for the FT-ICR-MS, the top ranking metabolite with this range indicated as the
192 identification for each discriminatory negative ionisation mode FIE-MS metabolite.

193

194 **2.8 Data and Statistical Analysis**

195 All GC-MS data pre-treatment procedures, including baseline correction, chromatogram alignment, and data
196 compression were performed by using custom scripts in Matlab version 6.5.1 (The Math Works Inc). Targeted
197 peak lists were generated, and peak apex intensities of each characteristic mass in a retention time window

198 were saved in an intensity matrix (run x metabolite). FIE-MS data was normalised with the total ion count for
199 each sample used to transform the intensity value for each metabolite in to a percentage of the total ion count,
200 after the removal of metabolites below 50 m/z. Principal Component Analyses (PCA). Briefly, this involves
201 projecting a ("X") matrix formed from set of (N) spectra with a given mass range (p) (i.e. a N x p matrix) onto
202 multi-dimensional space. Principal components (PC) are linear combinations of original variables (known as
203 loadings) which are used in the projection of the X matrix. Individual PCs are ranked (PC1, PC2...etc) on the basis
204 of the variance within the original data-set that is explained. PCA is an unsupervised method where no *a priori*
205 knowledge of experimental structure is given. Thus, if there is clustering of either 2D or 3D projections of PCA
206 from replicate data this indicates that the original experimental parameters are the sources of maximal
207 variation. Hierarchical Cluster Analyses (HCA) with heat maps, and Random Forest (RF) multivariate analysis
208 were completed using the PyChem (Version 3.0.5g Beta) package [13] and/or MetaboAnalyst 2.0 [14]. ROC
209 (Receiver Operating Characteristic) curve analyses plot the true positive rate (Sensitivity) in function of the false
210 positive rate (Specificity) and the validity of the fit is indicated based on area under curve (AUC) calculations.
211 ROC-AUC analyses used the ROC Curve Explorer and Tester (ROC CET) online platform [15] to assess our standard
212 binary classification tests. Due to the exploratory nature of this pilot study, no external validation set consisting
213 of an independent population of samples was available to be included in (e.g.) the ROC analyses.

214

215

216

217

218

219

220

221

222

223

224

225

226 3 Results

227

228 Patients and participants sampled as part of this study are summarised in Table 1, with individual sample data
229 in Supplementary Tables 1a and 1b. A total of 34 patients with suspected LC were recruited, with 23 confirmed
230 with LC (LC⁺) (16 NSCLC (nine Stage 4, three Stage 3A, one Stage 3B, three Stage 2B, and one Stage 1B), six SCLC
231 (three extensive and three limited), and one receiving a clinic-radiological diagnosis made by the LC
232 multidisciplinary team), and 11 had no diagnosis of LC after extensive testing and follow up for at least one year
233 (LC⁻). In addition, a total of 33 non-clinical controls (CON) were collected from participants with no history of
234 clinical lung disease.

235

236 Metabolomic profiles of the sputum samples were acquired using FIE-MS (in negative and positive ionization
237 modes) and GC-MS platforms and analysed using PCA. Both MS platforms were examined as although GC-MS is
238 widely employed in metabolomics profiling, it lacks the sensitivity of FIE-MS and thus, the latter could yield a
239 more comprehensive data set [16]. PCA indicated that the metabolomic profile acquired in negative ionisation
240 FIE-MS mode (Figure 1a) showed the greatest degree of separation between the three sample groups (LC⁺/LC⁻
241 /CON). Such a separation was not evident in positive FIE-MS mode (Figure 1b), and only partially exhibited in
242 the analyses of the GC-MS profiles (Figure 1c) suggestive of the value of the greater sensitivity of the FIE-MS
243 approach and platform. The FIE-MS⁻ metabolites were then analysed using one-way ANOVA which identified
244 the top 25 metabolites based on their discriminatory ability whose levels significantly differed between the
245 sample groups. Derivation of a HCA with heat map based on these top 25 metabolites also demonstrated that
246 the LC⁺ and LC⁻ could be readily separated from the CON group (Figure 2). Furthermore, many LC⁺ samples
247 clustered together.

248

249 Whilst simple analyses such as PCA or ANOVA could distinguish between the LC⁺/LC⁻ class and CON, supervised
250 analyses, where *a priori* information of the sample classes was required, would be need to identify variable
251 between the LC⁺ and LC⁻ classes. Due to the separation shown with FIE-MS⁻ metabolites into clinically relevant
252 classes these datasets were used to identify clinical relevant metabolomic biomarkers. Random Forest (RF)
253 analyses were then used to indicate a number of metabolites which differentiated between the experimental

254 classes (Figure 3). Metabolites which were either increased or decreased in the LC⁺ or CON classes compared to
255 the LC⁻ class which were taken forwards to identification by high resolution MS.

256

257 ROC -AUC analyses were also used to identify discriminatory metabolites. The top five metabolites for each
258 differential comparison are listed in Table 2 with the AUC figure and box and whisker distributions of the data
259 for the top differential metabolites shown in Figure 4. *t*-tests of the targeted metabolites indicated a high level
260 of significance in each comparison. These identified a number of metabolites that had a high AUC value (>0.99)
261 for differentiating between non-clinically (CON- class) and clinically acquired (LC/LC⁺ classes) samples. Four
262 metabolites were identified with an AUC value of greater than 0.80, a threshold for clinically useful prediction.

263

264 To identify the mass-ions targeted by RF and ROC-AUC analyses, high resolution MS using FT-ICR-MS was
265 employed. Metabolites were identified, where possible, based on this accurate mass profiling and database
266 interrogations, (Table 3) and these were used to annotate the analyses shown in Figures 3 and 4. Examination
267 of the metabolites listed in Supplementary Table 2 includes those involved in polyamine (putrescine), amino
268 acid, and lipid metabolism. Clinical samples (LC⁺/LC⁻) appeared to be separated from CON sample through
269 differential processing of polyamine metabolites; putrescine and N,N,N-Trimethylethenaminium, and lipid
270 metabolites, including glycerophospholipids of the cardiolipin (PC) class, and isobutyl decanoate and diethyl
271 glutarate. Separation between the clinical samples (LC⁺ and LC⁻ classes) appeared to be due to elevated levels of
272 metabolites identified as hexanal, cysteic acid, hydroxypyruvic acid, and the cholesterol ester with an acyl group
273 CE (22:5(4Z,7Z,10Z,13Z,16Z)). The mass-ion 1496.72 showing the highest AUC value (0.85) was identified as the
274 ganglioside GM1 (18:1/12:0).

275

276

277

278

279

280

281

282 4 Discussion

283

284 Since the 'Warburg Effect' was first described in 1956 [17], the alterations that cells undergo during
285 carcinogenesis has been a focus of both basic and applied clinical research. To date, the majority of metabolomic
286 lung cancer studies appear to have focussed on the cancerous tumours themselves or serum from affected
287 patients, using a limited range of MS techniques [18]. Here, we suggest that the sputum of patients can be used
288 as a non-invasive source of biomarkers for the identification of LC status.

289

290 Sputum represents a biofluid that could be readily accessed from the target group and the results of this study
291 indicate it could be used as a biofluid matrix for an efficient LC screen. We used two mass-spectrometry
292 platforms; the widely employed GC-MS and also FIE-MS on the same sample set to allow comparison of the
293 discriminatory power of both. These results suggested that derivatisation (in the case of GC-MS) or the wide
294 range of adducts formed with positive ionisation using FIE-MS (as opposed to negative, ionisation where simple
295 proton loss [$M^- - H^+$] is predominant) can obscure screens of sputum.

296

297 Analyses of FIE-MS⁻ data allowed identification of clinically relevant groupings and both PCA and HCA could
298 separate a "healthy" control samples from samples taken from clinically-referred patients. Although, not all of
299 these patients were subsequently confirmed to be LC⁺, the LC⁻ group had symptoms necessitating referral and
300 thus, should be considered to be "unwell". Even at this level, a non-invasive and rapid test of lung health would
301 be useful to the medical community.

302

303 Random Forest analysis appeared to be particularly effective in discriminating between LC⁺ and CON samples;
304 with LC⁻ samples between these extremes. We coupled RF analyses with assessment of ROC using AUC analysis;
305 which has been widely used to determine the diagnostic value of biomarkers. Here, the False Discovery rate vs.
306 True Discovery rate compared a series of binary tests between our three sample groups. The CON group was
307 highly distinctive, with ROC-AUC analyses detecting metabolites with extremely high AUC values. Crucially, a
308 number of FIE-MS⁻ metabolites that had AUC values greater than 0.80 when comparing LC⁺/LC⁻, a cut-off for
309 discrimination that may be useful in a clinical setting were identified. This equates to a false discovery rate of

310 under 20% for these metabolites although care needs to be taken with this figure and it requires confirmation
311 with external validation datasets to remove any danger of “overfitting” i.e. deriving a model which describes
312 random error or noise rather than any true relationship. Our LC⁺ group consisted of a range of LC stages and
313 histology, suggesting that biomarkers established through metabolomic profiling techniques could have utility
314 as a preliminary screen, identifying patients for clinical follow-up for LC confirmation, histology and staging.

315

316 Considering the identities of metabolites separating the CON and LC⁻ class, the increases in putrescine were
317 interesting because polyamines are essential for normal mammalian cell growth. Polyamine metabolism is
318 frequently dysregulated in cancer and has emerged as a target for therapeutic intervention [19]. However, as
319 polyamines did not discriminate between the CON and LC⁺ or LC⁻ and LC⁺ classes, we were not able to associate
320 these polyamine changes with LC in this study. Therefore, changes in polyamines may have reflected changes
321 linked to an inflammatory response and/or cell death; which may reflect pathogen attack or polyamine
322 catabolism which can generate reactive oxygen species (ROS) [20].

323

324 Also prominent in the clinical samples (LC⁺/LC⁻) compared to the CON class were lipid metabolites, including
325 glycerophospholipids of the cardiolipin (PC) class as well as isobutyl decanoate and diethyl glutarate. Cardiolipins
326 are major components of the inner mitochondrial and is particularly susceptible to ROS attack due to its high
327 content of unsaturated fatty acids. Increased ROS would affect mitochondrial membrane fluidity, possibly
328 resulting in cardiolipin release and possibly leading to the greater than two fold increases that we have detected
329 in our study. Cardiolipin-associated changes in membrane fluidity have been associated with reduced
330 mitochondrial oxidative phosphorylation efficiency and apoptosis [21]. In this context, it is relevant that isobutyl
331 decanoate and diethyl glutarate, as potential phospholipid fragments, could represent the products of lipid
332 peroxidation and as they exhibited a 4.69 and 3 fold increase, respectively, in the LC⁺ class compared to CON.

333

334 Identifying the metabolite changes in the LC⁺ samples targeted by RF and ROC-AUC, there appeared to be higher
335 levels of hexanal, cysteic acid, hydroxypyruvic acid, and one metabolite without accurate mass identification,
336 and eleven metabolites with lower levels. Hexanal has previously been shown to be elevated in blood samples
337 from lung cancer patients [22], suggesting its validity as a LC biomarker. To our knowledge there have been no

338 reports of cysteic acid or hydroxypyruvic acid being targeted as LC biomarkers. These could suggest alterations
339 in cysteine metabolism (in the case of cysteate) or glycolysis (in the case of the pyruvic acid derivative) are being
340 targeted in our metabolomic analyses. The potential relevance of the cholesteryl docosapentaenoate;
341 CE(22:5(4Z,7Z,10Z,13Z,16Z)) is unknown, but its increase could reflect membrane disruption. Of particular
342 interest was ganglioside GM1 (18:1/12:0) which represents a glycosphingolipid inked to a single sialic acid
343 through its sugar group. Gangliosides have primarily been studied in neural tissues, but can be found in most
344 cell types where they are involved in cell–cell recognition, cell–matrix attachment, cell growth and cell
345 differentiation. Interestingly, ganglioside GM1 has already been associated with LC and particularly with SCLC
346 due to a tendency to arise from neuroectodermal tissue [23]. Indeed, GM1 ganglioside-fused to hemocyanin
347 has been used to specifically target SCLC tissue in patients [24]. Cholera toxin which is known to target GM1
348 ganglioside was found to suppress the growth of 9 out of 15 SCLC cell lines with those resistant to the toxin
349 exhibiting reduced GM1 ganglioside expression [25]. Taken together with our results, ganglioside GM1 could be
350 a good candidate for biomarker based LC screens.

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366 **5 Conclusions**

367

368 As far as we can ascertain, this is the first study to report on the metabolomic profiling of sputum acquired from
369 LC patients. The use of sputum, the production of which is symptomatic of LC, as a biofluid for screening carries
370 the benefit of being non-invasive, high-throughput, and low-cost, compared to current conventional methods
371 such as CT scan [26]. It may be that a combination of metabolomic biomarkers and other types, such as
372 circulating miRNAs, would allow for an integrated approach to LC screening, as has been suggested for other
373 cancers [27]. Here, we have shown the power of using metabolomics to identify biomarkers with potential
374 clinical application for LC. Further work, using a larger patient cohort, will be required to ascertain the utility of
375 metabolomic biomarkers for LC stage and histological subtype.

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394 **Acknowledgements:** SJSC is grateful for studentship from Aberystwyth University. The MedLung Study is funded
395 by a grant from the National Institute of Social Care and Health Research (NISCHR), Wales. The sponsor was
396 Hywel Dda University Health Board and neither the funders – Aberystwyth University or NISCHR - nor sponsor
397 had any input into the design or reporting of the study. We wish to thank Dr Paul Griffiths, Consultant
398 histopathologist, for sputum samples assessment, Dr Sion Bayliss for collection of healthy age-matched control
399 samples, and Kathleen Taillart for running samples through mass spectrometry. We are also highly appreciative
400 of the constructive criticisms provided by the anonymous reviewers which improved this manuscript.

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422 **References**

423

- 424 [1] WHO, Cancer Factsheet, WHO Fact Sheets, Number 2970. (2013).
425 <http://www.who.int/mediacentre/factsheets/fs297/en/> (accessed May 29, 2013).
- 426 [2] A. Jemal, R. Siegel, J.Q. Xu, E. Ward, Cancer Statistics, 2010, *CA A Cancer J. Clin.* 60 (2010) 277–300.
427 doi:10.1002/caac.20073.
- 428 [3] G. Sutedja, New Techniques for Early Detection of Lung Cancer., *Eur. Respir. J.* 39 (2003) 57s–66s.
429 <http://www.ncbi.nlm.nih.gov/pubmed/12572703> (accessed February 25, 2014).
- 430 [4] P. Jha, R. Peto, W. Zatonski, J. Boreham, M.J. Jarvis, A.D. Lopez, Social Inequalities in Male Mortality, and
431 in Male Mortality from Smoking: Indirect Estimation from National Death Rates in England and Wales,
432 Poland, and North America., *Lancet.* 368 (2006) 367–70. doi:10.1016/S0140-6736(06)68975-7.
- 433 [5] Y. Xie, N.W. Todd, Z. Liu, M. Zhan, H. Fang, H. Peng, et al., Altered miRNA Expression in Sputum for
434 Diagnosis of Non-Small Cell Lung Cancer., *Lung Cancer.* 67 (2010) 170–6.
435 doi:10.1016/j.lungcan.2009.04.004.
- 436 [6] J. Shen, N.W. Todd, H. Zhang, L. Yu, X. Lingxiao, Y. Mei, et al., Plasma microRNAs as Potential Biomarkers
437 for Non-Small-Cell Lung Cancer., *Lab. Investig.* 91 (2011) 579–87. doi:10.1038/labinvest.2010.194.
- 438 [7] K.M. Foss, C. Sima, D. Ugolini, M. Neri, K.E. Allen, G.J. Weiss, miR-1254 and miR-574-5p: Serum-Based
439 microRNA Biomarkers for Early-Stage Non-Small Cell Lung Cancer., *J. Thorac. Oncol.* 6 (2011) 482–8.
440 doi:10.1097/JTO.0b013e318208c785.
- 441 [8] P.D. Lewis, K.E. Lewis, R. Ghosal, S. Bayliss, A.J. Lloyd, J. Wills, et al., Evaluation of FTIR Spectroscopy as a
442 Diagnostic Tool for Lung Cancer Using Sputum, *BMC Cancer.* 10 (2010). doi:10.1186/1471-2407-10-
443 640.
- 444 [9] J. Draper, A.J. Lloyd, R. Goodacre, M. Beckmann, Flow infusion electrospray ionisation mass spectrometry
445 for high throughput, non-targeted metabolite fingerprinting: a review, *Metabolomics.* 9 (2012) 4–29.
446 doi:10.1007/s11306-012-0449-x.
- 447 [10] G.S. Collins, J.B. Reitsma, D.G. Altman, K.G.M. Moons, Transparent reporting of a multivariable prediction
448 model for individual prognosis or diagnosis (TRIPOD): the TRIPOD Statement., *BMC Med.* 13 (2015) 1.
449 doi:10.1186/s12916-014-0241-z.

- 450 [11] A.J. Lloyd, G. Favé, M. Beckmann, W. Lin, K. Taillart, L. Xie, et al., Use of Mass Spectrometry Fingerprinting
451 to Identify Urinary Metabolites after Consumption of Specific Foods., *Am. J. Clin. Nutr.* 94 (2011) 981–91.
452 doi:10.3945/ajcn.111.017921.
- 453 [12] D.S. Wishart, T. Jewison, A.C. Guo, M. Wilson, C. Knox, Y. Liu, et al., HMDB 3.0: The Human Metabolome
454 Database in 2013., *Nucleic Acids Res.* 41 (2013) D801–7. doi:10.1093/nar/gks1065.
- 455 [13] R.M. Jarvis, D. Broadhurst, H. Johnson, N.M. O’Boyle, R. Goodacre, PYCHEM: A Multivariate Analysis
456 Package for Python., *Bioinformatics.* 22 (2006) 2565–6. doi:10.1093/bioinformatics/btl416.
- 457 [14] J. Xia, R. Mandal, I. V Sineelnikov, D. Broadhurst, D.S. Wishart, MetaboAnalyst 2.0: A Comprehensive Server
458 for Metabolomic Data Analysis., *Nucleic Acids Res.* 40 (2012) W127–33. doi:10.1093/nar/gks374.
- 459 [15] J. Xia, D.I. Broadhurst, M. Wilson, D.S. Wishart, Translational Biomarker Discovery in Clinical
460 Metabolomics: An Introductory Tutorial., *Metabolomics.* 9 (2013) 280–299. doi:10.1007/s11306-012-
461 0482-9.
- 462 [16] L.W. Sumner, P. Mendes, R.A. Dixon, Plant metabolomics: large-scale phytochemistry in the functional
463 genomics era, *Phytochemistry.* 62 (2003) 817–836. doi:10.1016/S0031-9422(02)00708-2.
- 464 [17] O. Warburg, On the Origin of Cancer Cells, *Science* (80-.). 123 (1956) 309–314.
465 doi:10.1126/science.123.3191.309.
- 466 [18] S. Hori, S. Nishiumi, K. Kobayashi, M. Shinohara, Y. Hatakeyama, Y. Kotani, et al., A Metabolomic
467 Approach to Lung Cancer., *Lung Cancer.* 74 (2011) 284–92. doi:10.1016/j.lungcan.2011.02.008.
- 468 [19] A.E. Pegg, Mammalian Polyamine Metabolism and Function., *IUBMB Life.* 61 (2009) 880–94.
469 doi:10.1002/iub.230.
- 470 [20] M.H. Park, K. Igarashi, Polyamines and their Metabolites as Diagnostic Markers of Human Diseases.,
471 *Biomol. Ther. (Seoul).* 21 (2013) 1–9. doi:10.4062/biomolther.2012.097.
- 472 [21] G. Paradies, G. Petrosillo, V. Paradies, F.M. Ruggiero, Role of Cardiolipin Peroxidation and Ca²⁺ in
473 Mitochondrial Dysfunction and Disease., *Cell Calcium.* 45 (2009) 643–50.
474 doi:10.1016/j.ceca.2009.03.012.
- 475 [22] C. Deng, X. Zhang, N. Li, Investigation of volatile biomarkers in lung cancer blood using solid-phase
476 microextraction and capillary gas chromatography-mass spectrometry., *J. Chromatogr. B. Analyt.*
477 *Technol. Biomed. Life Sci.* 808 (2004) 269–77. doi:10.1016/j.jchromb.2004.05.015.

- 478 [23] T. Brezicka, B. Bergman, S. Olling, P. Fredman, Reactivity of Monoclonal Antibodies with Ganglioside
479 Antigens in Human Small Cell Lung Cancer Tissues, *Lung Cancer*. 28 (2000) 29–36. doi:10.1016/S0169-
480 5002(99)00107-5.
- 481 [24] L.M. Krug, G. Ragupathi, C. Hood, M.G. Kris, V.A. Miller, J.R. Allen, et al., Vaccination of patients with
482 small-cell lung cancer with synthetic fucosyl GM-1 conjugated to keyhole limpet hemocyanin., *Clin.*
483 *Cancer Res*. 10 (2004) 6094–100. doi:10.1158/1078-0432.CCR-04-0482.
- 484 [25] R. Fuentes, R. Allman, M.. Mason, Ganglioside expression in lung cancer cell lines, *Lung Cancer*. 18 (1997)
485 21–33. doi:10.1016/S0169-5002(97)00049-4.
- 486 [26] M.P. Rivera, A.C. Mehta, M.M. Wahidi, Establishing the Diagnosis of Lung Cancer: Diagnosis and
487 Management of Lung Cancer, *Chest*. 143 (2013) e142S–65S. doi:10.1378/chest.12-2353.
- 488 [27] B. Laxman, D.S. Morris, J. Yu, J. Siddiqui, J. Cao, R. Mehra, et al., A First-Generation Multiplex Biomarker
489 Analysis of Urine for the Early Detection of Prostate Cancer., *Cancer Res*. 68 (2008) 645–9.
490 doi:10.1158/0008-5472.CAN-07-3224.

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506 **Tables and Legends**

507

508 **Table 1 | Summarised Patient and Participant Information**

509 Summarised patient information detailing clinical data. Full clinical data for clinically acquired samples, and
 510 information collected for healthy control participants, are fully detailed in Supplementary Table 1.

511

	Non-Clinical Controls (CON)	LC Negative (LC-)	LC Positive (LC+)
Number	33	11	23
Age	55.3 (14.6)	66.5 (14.3)	66.6 (8.1)
Gender			
Male	20	10	11
Female	13	1	12
Smoking Status			
Current	15	3	10
Ex	0	8	10
Never	18	0	3
Smoking Pack Years	NC	49.0 (34.9)	39.3 (18.9)
Infection Present			
Yes	NC	3	1
No	NC	8	22
CO Level (ppm)	NC	3.7 (1.3)	4.2 (2.8)

512

513

514

515

516

517

518

519

520

521

522

523 **Table 2 | Top Five Area Under Curve Values for Negative FIE-MS Mode Metabolites**

524 Using the online ROC CET platform, the top five metabolites, based on AUC values, for each differential group
 525 comparison were identified. For clinical and non-clinical comparisons, high AUC values were obtained, and for
 526 the LC negative and positive comparison, a number of metabolites were identified with AUC values greater than
 527 0.8. AUC range refers to the 95% confidence intervals of the true AUC value as given by the ROC CET platform.

528

Differential	Metabolite	AUC Value	True AUC Range	t-Test	Fold Change
CON Vs LC-	N,N,N-Trimethylethenaminium /CL(16:1(9Z)/18:1(11Z)/16:1(9Z)/18:1(9Z))	1.00	0.989 - 1.000	4.47×10^{-15}	-2.38
	N,N,N-Trimethylethenaminium /1560.81	1.00	0.983 - 1.000	3.10×10^{-15}	-2.26
	Putrescine/CL(16:1(9Z)/18:1(11Z)/16:1(9Z)/18:1(9Z))	0.99	0.975 - 1.000	6.09×10^{-15}	-2.25
	Putrescine/1560.81	0.99	0.975 - 1.000	3.64×10^{-15}	-2.14
	53.27/1209.45	0.99	0.967 - 1.000	1.79×10^{-14}	-2.28
LC+ Vs CON	53.27/Isobutyl decanoate	1.00	0.993 - 1.000	8.74×10^{-24}	4.71
	Putrescine/Isobutyl decanoate	1.00	0.992 - 1.000	7.42×10^{-24}	4.69
	189.09	1.00	0.987 - 1.000	6.12×10^{-20}	2.57
	Diethyl glutarate	0.99	0.979 - 1.000	1.42×10^{-18}	3.00
	Cysteamine	0.99	0.980 - 1.000	3.03×10^{-20}	2.35
LC+ Vs LC-	Ganglioside GM1 (18:1/12:0)	0.85	0.709 - 0.953	2.93×10^{-3}	-0.03
	957.36	0.83	0.680 - 0.953	4.57×10^{-3}	0.31
	1382.45	0.83	0.668 - 0.957	5.93×10^{-4}	0.07
	CE(22:5(4Z,7Z,10Z,13Z,16Z))	0.82	0.644 - 0.947	1.28×10^{-2}	0.00
	1434.00	0.81	0.621 - 0.947	9.84×10^{-4}	0.14

529

530

531 **Figures**

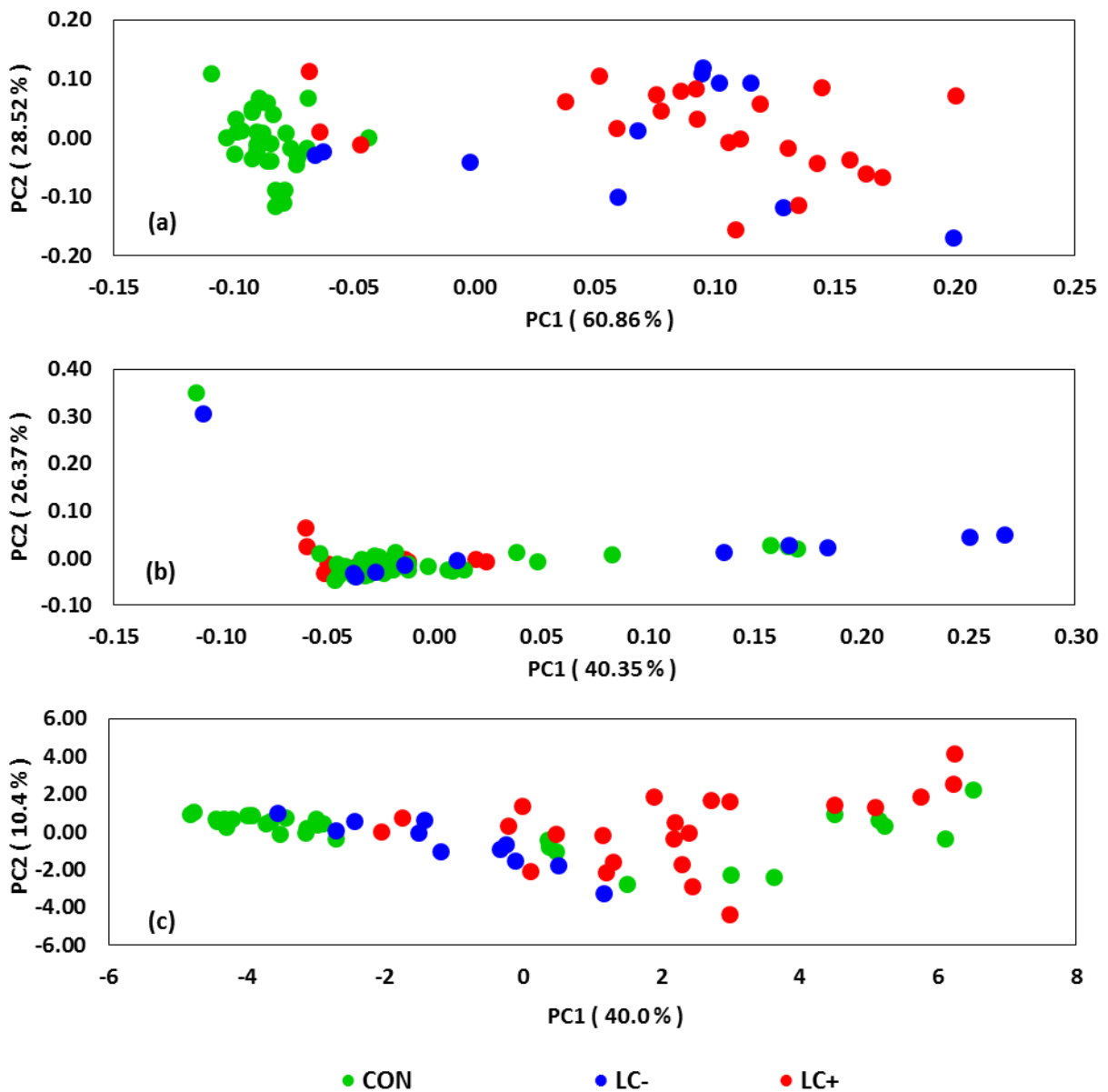
532

533 **Figure 1 | Principal Component Analysis Plots for FIE-MS and GC-MS Metabolites**

534 PCA, based on metabolites acquired in (a) FIE-MS negative mode, (b) FIE-MS positive mode, and (c) GC-MS,
535 clearly differentiates between the clinically and non-clinically acquired samples, though separation of the two
536 clinical groups, lung cancer and symptom controls, does not occur. For (c), coordinate markers are means of
537 individually calculated coordinates from duplicate GC-MS runs.

538

539



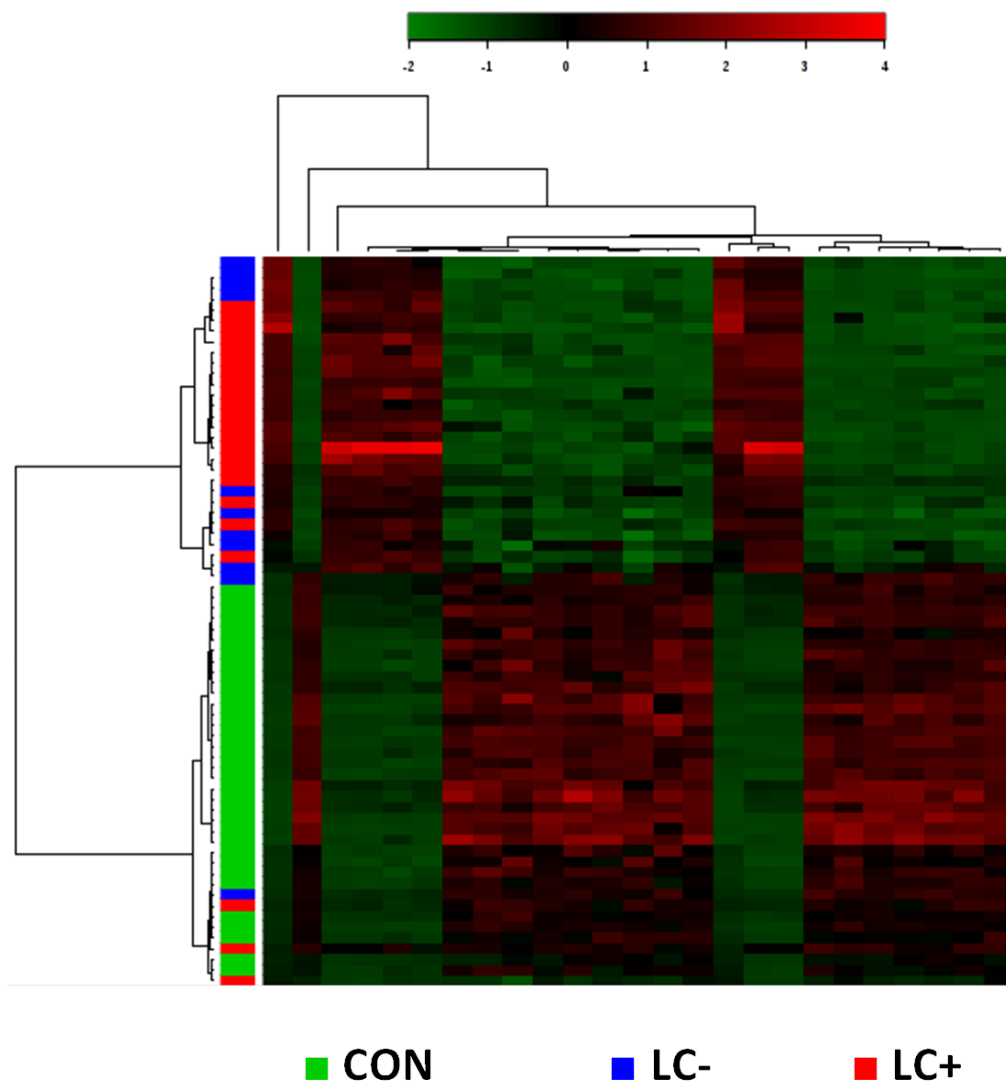
540

541

542 **Figure 2 | Hierarchical Cluster Analysis with Heat Mapping for Negative Ionisation FIE-MS**

543 Hierarchical cluster analysis and corresponding heat maps were constructed, based on the top 25 metabolites
544 identified through one-way ANOVAs, for metabolites identified in FIE-MS negative ionisation mode. Similarly to
545 PCA plots, separation between the clinically and non-clinically acquired samples was clear, but separation
546 between LC positive and negative samples was not evident.

547



548

549

550

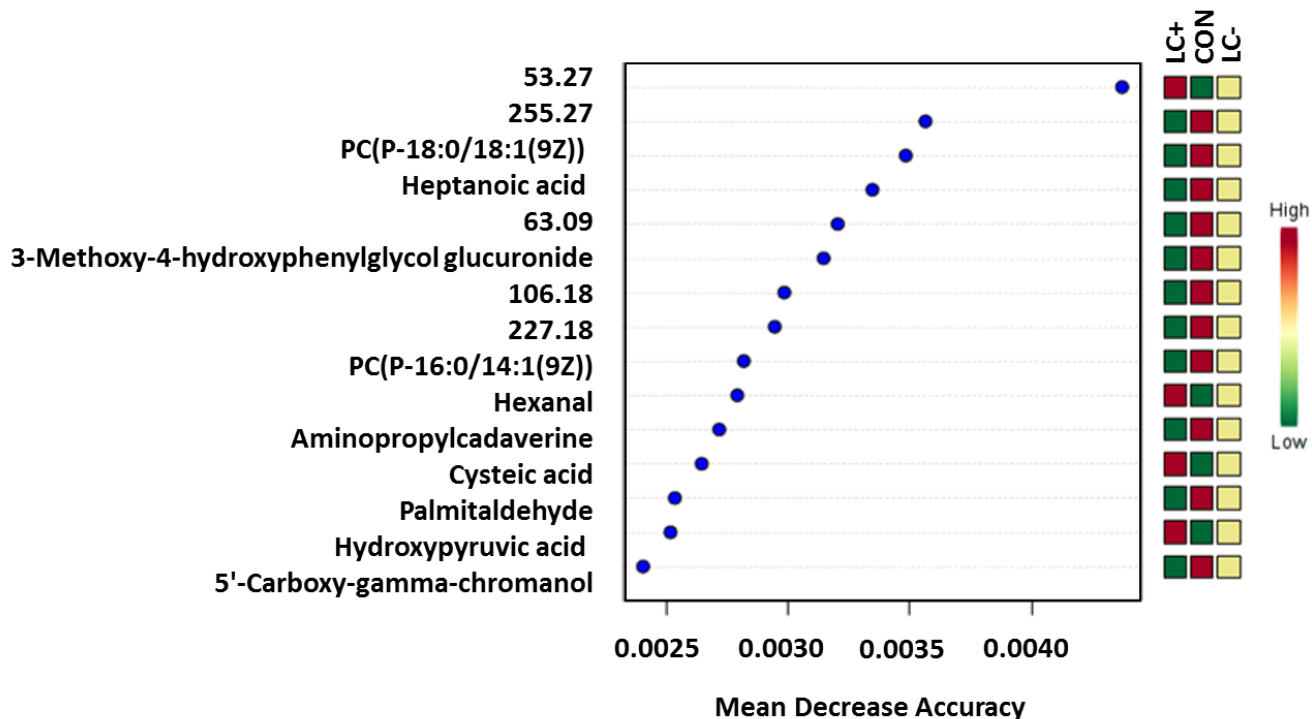
551

552

553

554 **Figure 3 | Random Forest Plots for Identification of Key FIE-MS Metabolites**

555 Random forests plots were constructed, using MetaboAnalyst 2.0 for negative ionisation FIE-MS mode, which
 556 revealed a number of metabolites which may have potential in terms of diagnostic markers, particularly those
 557 that are either higher or lower in the LC positive group.

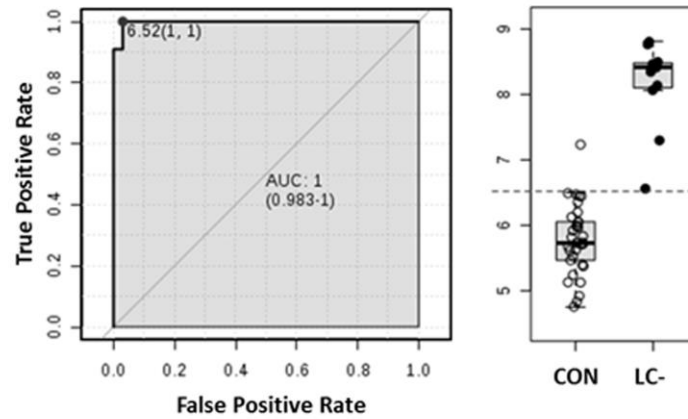


558
 559
 560
 561
 562
 563
 564
 565
 566
 567
 568
 569
 570

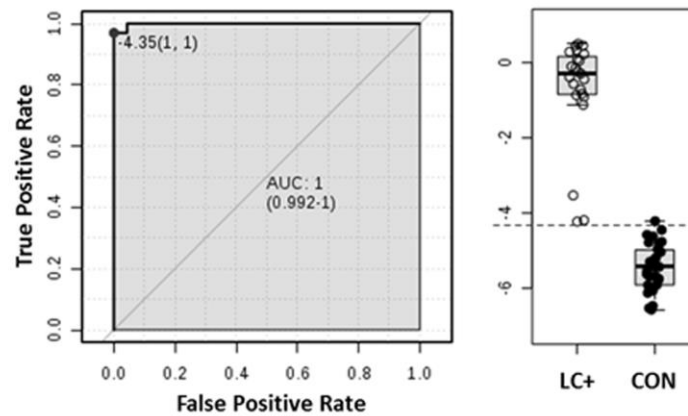
571 **Figure 4 | Univariate Receiver Operating Characteristic Curve Analyses for Biomarker Identification**

572 Using the online facility, ROCET, univariate receiver operating characteristic curves (ROC) were created, and
573 plotted to create area under the curve (AUC) figures for metabolites identified in negative ionisation FIE-MS
574 mode. The metabolite with the highest AUC value for each differential group is plotted.

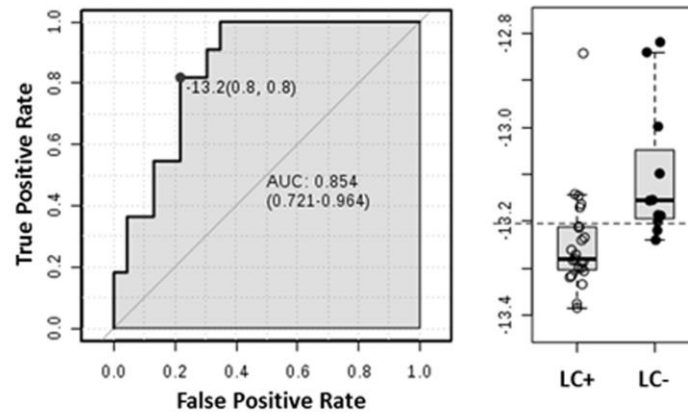
**(a) N,N,N-Trimethylethenaminium /
CL(16:1(9Z)/18:1(11Z)/16:1(9Z)/18:1(9Z))**



(b) 53.27 / Isobutyl decanoate



(c) Ganglioside GM1 (18:1/12:0)



575

576