

This document is confidential and is proprietary to the American Chemical Society and its authors. Do not copy or disclose without written permission. If you have received this item in error, notify the sender and delete all copies.

**A machine learning approach using a handheld near-infrared (NIR) device to predict the effect of storage conditions on tomato biomarkers**

Journal:	<i>ACS Food Science &amp; Technology</i>
Manuscript ID	Draft
Manuscript Type:	Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Emsley, Natalia; Lancaster Girls' Grammar School, Chemistry Holden, Claire; Lancaster University Faculty of Science and Technology Guo, Yi; Lancaster Girls' Grammar School, Chemistry Bevan, Rhiann; Lancaster Girls' Grammar School, Chemistry Rees, Chris; Lancaster Girls' Grammar School, Chemistry McAinsh, Martin; Lancaster University, Centre for Biophotonics Martin, Francis; Biocel Ltd, Biocel Analytics Morais, Camilo; University of Central Lancashire, School of Pharmacy and Biomedical Sciences

SCHOLARONE™  
Manuscripts

# A machine learning approach using a handheld near-infrared (NIR) device to predict the effect of storage conditions on tomato biomarkers

Natalia E M Emsley<sup>1,¶</sup>, Claire A. Holden<sup>2,¶</sup>, Sarah Guo<sup>1</sup>, Rhiann S Bevan<sup>1</sup>, Christopher Rees<sup>3</sup>, Martin R. McAinsh<sup>2</sup>, Francis L Martin<sup>4,\*</sup>, Camilo L M Morais<sup>5,\*</sup>

*Sixth Form, Lancaster Girls' Grammar School, Lancaster, Lancashire, UK*

*<sup>2</sup> Lancaster Environment Centre, Library Avenue, Lancaster University, Lancaster LA1 4YQ, UK*

*<sup>3</sup> Chemistry Department, Lancaster Girls' Grammar School, Lancaster, Lancashire, UK*

*<sup>4</sup> Biocel Ltd, Hull HU10 7TS, UK*

*<sup>5</sup> School of Pharmacy and Biomedical Sciences, University of Central Lancashire, Preston, Lancashire, UK*

**\*Corresponding authors:** Francis L Martin ([film13@biocel.uk](mailto:film13@biocel.uk)); Camilo L M Morais ([camilomorais1@gmail.com](mailto:camilomorais1@gmail.com))

¶, Joint first authors

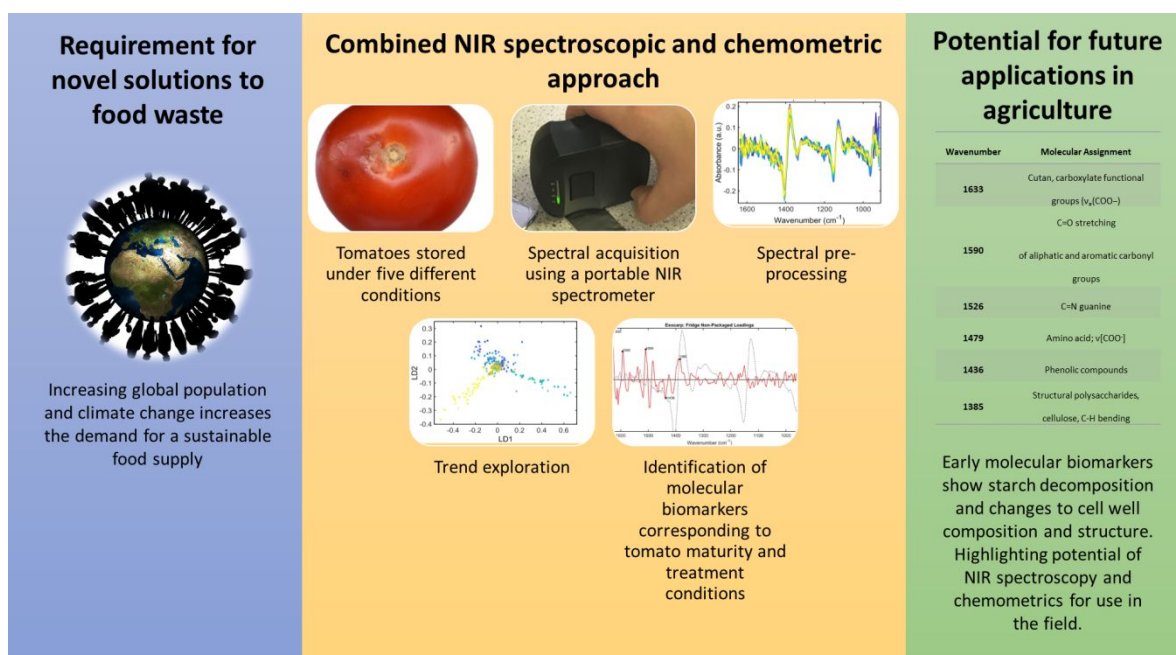
## 21 Abstract

22 Minimising food waste critical to future global food security. This study aimed to assess the potential  
 23 of near-infrared (NIR) spectroscopy combined with machine learning to monitor the stability of  
 24 tomato fruit during storage. Freshly harvested UK-grown tomatoes ( $n=135$ ) were divided into five  
 25 equally sized groups, each stored in different conditions. Absorbance spectra were obtained from  
 26 both the tomato exocarp and locular gel using a portable NIR spectrometer, capable of connecting  
 27 to a mobile phone, before subsequent chemometric analysis. Results show that support vector  
 28 machines can predict the storage conditions and time-after-harvest of tomatoes. Molecular  
 29 biomarkers highlighting key wavenumber and molecular changes due to time and storage conditions  
 30 were also identified. This method shows potential for development of this approach for use in the  
 31 field to help mitigate the environmental and economic impacts of food waste.

## 32 Keywords

33 Tomato, infrared spectroscopy, food security, machine learning, chemometrics, food storage

## 35 Table of Contents Graphic



## 37 Introduction

38 Sustainable food supply chains are required to provide adequate nutrition for the rising global  
39 population, which is projected to reach approximately 10 billion by 2050 <sup>1</sup>. Achieving food security  
40 and ending malnutrition are priorities in the United Nations sustainable development agenda for  
41 2030 <sup>2</sup>. To achieve these goals using the same amount of available land, it is important that to  
42 maximise the efficiency of current supply chains; this includes the reduction of food waste. In the  
43 UK, food waste within the home is largely avoidable with 2.9 million metric tons, equating to £6.3  
44 billion, of food spoiling before consumption <sup>3</sup>. Tomatoes are a popular crop globally, representing a  
45 versatile and nutrient-dense superfood, offering an excellent source of lycopene and other  
46 antioxidant molecules which benefit human health <sup>4</sup>. In 2017, 9.2 kg of fresh tomatoes and 33.2 kg  
47 of processed tomato products were consumed per capita in the United States of America <sup>5</sup>.  
48 However, fresh ripe tomatoes are prone to a high rate of food waste, with a total loss of 53.8%  
49 collectively from production to consumption <sup>6</sup>. This results from a 20% loss in agricultural  
50 production, 7% loss in processing and packaging, 10% distribution and retail loss and 31% loss by  
51 consumers <sup>6</sup>. The importance of reducing this food waste is emphasised during the current climate  
52 crisis, particularly as the production of tomatoes comes with a high environmental impact. They  
53 have been recognized as one of the most carbon-intensive food products due to electricity and  
54 fertiliser usage <sup>7</sup>, and greenhouse grown tomatoes are estimated to have a carbon footprint of up to  
55 10.1 kg CO<sub>2</sub>-eq/kg <sup>8</sup>. Improving the efficiency of the supply chain for tomato, and agri-food in  
56 general, would reduce the economic and environmental impacts of food waste, and help us to meet  
57 the challenge of nourishing a growing population. Technical innovations to increase shelf-life have  
58 the potential to reduce the unnecessary disposal of home produce and hence reduce food waste <sup>9</sup>.  
59 Factors that affect the marketability and shelf-life of fresh ripe tomatoes include the duration that  
60 the fruits are safe for consumption, but also aesthetic qualities such as colour and firmness. As a  
61 climacteric fruit, the quality of fresh tomatoes can be improved after harvest. The ripening process  
62 involves changes to cell wall composition and thickness <sup>10</sup>, and the conversion of storage

1  
2  
3 63 carbohydrates into monosaccharides <sup>11</sup>. These processes increase the softness and sweetness of the  
4  
5 64 fruit, which in turn increases its appeal to consumers. However, if damaged or left to ripen for too  
6  
7 65 long, tomatoes are prone to infections such as sour-rot <sup>12</sup>. Early biomarker determination to rapidly  
8  
9 66 predict fruit ripening, and hence its potential shelf-life before the fruit are no longer saleable, could  
10  
11 67 help to achieve the optimum balance in post-harvest ripening under different climacteric conditions.  
12  
13  
14  
15 68 Near-infrared (NIR; wavelengths from 800 to 2,500 nm) spectroscopy has been used extensively in  
16  
17 69 various recent studies to quantitatively determine fruit quality or maturity in tomatoes <sup>12-18</sup>. The  
18  
19 70 portability of new miniaturised spectrometers has opened up new applications for spectroscopic  
20  
21 71 studies, such as agriculture, see <sup>19-22</sup>. Near infrared spectroscopy uses changes in molecular  
22  
23 72 vibrations upon absorption of infrared light to gain information about the chemical composition of a  
24  
25 73 sample. Biological materials preferentially absorb light in the fingerprint region (1800-900  $\text{cm}^{-1}$ ),  
26  
27 74 which gives information about key biomolecules <sup>23</sup>. These include: lipids (C=O symmetric stretching  
28  
29 75 at  $\sim 1,750 \text{ cm}^{-1}$  and CH<sub>2</sub> bending at  $\sim 1,470 \text{ cm}^{-1}$ ), proteins (amide I at  $\sim 1,650 \text{ cm}^{-1}$ , amide II at  
30  
31 76  $\sim 1,550 \text{ cm}^{-1}$  and amide III at  $\sim 1,260 \text{ cm}^{-1}$ ), carbohydrates (CO-O-C symmetric stretching at  $\sim 1,155$   
32  
33 77  $\text{cm}^{-1}$ ), nucleic acid (asymmetric phosphate stretching at  $\sim 1,225 \text{ cm}^{-1}$  and symmetric phosphate  
34  
35 78 stretching at  $\sim 1,080 \text{ cm}^{-1}$ ), glycogen (C-O stretching at  $\sim 1,030 \text{ cm}^{-1}$ ) and protein phosphorylation  
36  
37 79 ( $\sim 970 \text{ cm}^{-1}$ ) <sup>23</sup>. The current study utilised a handheld NIR spectroscopy device to determine tomato  
38  
39 80 maturity and explore the effects of five different storage conditions over twenty days, varying in  
40  
41 81 temperature and packaging. Spectral measurements were taken from both the exocarp and the  
42  
43 82 locular gel within, allowing comparison of results from destructive and conservative techniques. A  
44  
45 83 novel approach combining NIR using a handheld spectrometer and machine learning was used for  
46  
47 84 classification and identification of key biomarkers. Chemometric methods included principal  
48  
49 85 component analysis (PCA), PCA coupled with linear discriminant analysis (PCA-LDA) and support  
50  
51 86 vector machines (SVM). Our results highlight key wavenumber changes associated with ripening and  
52  
53 87 the effects of post-harvest climatic conditions, raising opportunity of developing this combined  
54  
55 88 method for use in the field.  
56  
57  
58  
59  
60

## 89 **Materials and Methods**

### 90 **Samples and data acquisition**

91 This study measured tomatoes over a period of twenty days under five different storage conditions  
92 (ambient packaged, ambient non-packaged, fridge packaged, fridge non-packaged, and incubated  
93 non-packaged). A total of 135 tomatoes were used in this study, twenty-seven per storage condition.  
94 All tomatoes were F1 hybrids, a first generation cross between two varieties, grafted onto a strong  
95 disease resistant rootstock and grown in Rockwool, UK (supplied by John Lane). They were harvested  
96 at ten weeks old on Friday 13<sup>th</sup> September 2019. The tomatoes used in this study were selected for  
97 uniformity. At the start of the study, they were all graded to be of similar size, colour, and ripeness,  
98 were undamaged, and had a consistent colour across the whole surface. They were then randomly  
99 sorted into treatment groups. Three tomatoes from each treatment condition were selected for  
100 exocarp sampling and labelled X, Y and Z. Each treatment group was stored in different conditions  
101 over the course of the experiment, including ambient packaged, ambient non-packaged, fridge  
102 packaged, fridge non-packaged and incubated non-packaged. The packaging used consisted of a  
103 sealed plastic freezer bag to simulate modified atmosphere packaging (MAP) although without the  
104 introduction of a protective gas mix into the bag. The ambient tomatoes were stored at room  
105 temperature (18°C), the fridge tomatoes at 3°C and the incubated tomatoes at 25°C.

106 Spectral absorbances were measured using a hand-held NIR spectrometer NIR-S-G1 (Allied Scientific  
107 Pro, Gatineau, Quebec, Canada) using ISC NIRScan software (raw spectra are shown in the  
108 Supporting Information, Figure S1). Spectra were measured from two types of sample; exocarp (the  
109 tomato surface), and locular gel (a gel that develops prior to ripening of the pericarp and exhibits a  
110 liquid-like consistency towards the terminal stage of ripening, see <sup>24</sup>). The spectrometer crystal was  
111 cleaned between measurements using isopropyl alcohol wipes (Bruker Optics, Coventry, UK), and  
112 each time background spectra were taken to account for ambient atmospheric conditions. For  
113 locular gel samples spectral acquisition took place on days 3, 5, 7, 10, 12, 14, 17, and 20. Each day  
114 three different tomatoes were destructively analysed to extract the locular gel, and two spectral

1  
2  
3 115 replicates per tomato were taken. This resulted in six locular gel spectra per day, forty-eight per  
4  
5 116 treatment group, and 240 spectra in total. The locular gel was extracted by incision at two  
6  
7 117 approximately equidistant locations around the equator. Fifty  $\mu\text{L}$  of this gel was collected using a  
8  
9 118 Gilson pipette and transferred onto glass slides covered in aluminium foil. Before spectral acquisition  
10  
11 119 of samples an aluminium foil standard spectrum was measured. For the exocarp measurements, the  
12  
13 120 same three tomatoes were analysed throughout the whole time-course and subsequently returned  
14  
15 121 to their original storage conditions. For exocarp samples, measurements were taken on days 1, 3, 5,  
16  
17 122 7, 10, 12, 14, 17, and 20. Ten spectral replicates of each tomato were taken, resulting 270 spectra  
18  
19 123 per treatment group and 1350 spectra in total.  
20  
21  
22

#### 23 124 [Data analysis and validation](#)

25  
26 125 The reflectance spectral data were imported and processed within MATLAB R2014b (MathWorks,  
27  
28 126 Inc., Natick, MA, USA). Pre-processing and data analysis were performed using the PLS Toolbox  
29  
30 127 version 7.9.3 (Eigenvector Research, Inc., Manson, WA, USA). The raw spectra were pre-processed  
31  
32 128 by Savitzky-Golay smoothing (window of 7 points, 2<sup>nd</sup> order polynomial fitting) to improve the signal-  
33  
34 129 to-noise ratio and standard normal variate (SNV) to correct for light scattering. The pre-processed  
35  
36 130 data were also mean-centred before multivariate analysis.  
37  
38  
39

40 131 Principal component analysis (PCA) was applied to the pre-processed spectral data for exploratory  
41  
42 132 analysis. PCA is an exploratory analysis technique that reduces the spectral dataset into a small  
43  
44 133 number of principal components (PCs) responsible for the majority of the original data variance <sup>25</sup>.  
45  
46 134 Each PC is orthogonal to each other, and they are formed in a decreasing order of explained  
47  
48 135 variance, so PC1 covers more variance than PC2, and so on. Each PC is composed of scores and  
49  
50 136 loadings. The scores represent the variance on sample direction, thus being used to identify patterns  
51  
52 137 of similarity between the samples, and the loadings represent the variance on wavelength direction,  
53  
54 138 thus being used to identify possible spectral markers responsible for the scores pattern. Although  
55  
56 139 being a robust exploratory analysis technique, PCA was not able to classify samples in an objective  
57  
58 140 fashion, therefore, a supervised classifier was added and the samples were analysed by principal  
59  
60

1  
2  
3 141 component analysis with linear discriminant analysis <sup>26</sup>. For more rigorous classification, support  
4  
5 142 vector machines (SVM) algorithm was applied to estimate the tomatoes time-after-harvest. SVM is a  
6  
7 143 linear classifier with a non-linear step called the kernel transformation <sup>27</sup>. The kernel function  
8  
9  
10 144 transforms the data space to a feature space where samples can be better discriminated. Herein, the  
11  
12 145 radial basis function (RBF) kernel was used and optimised via cross-validation. The pre-processed  
13  
14 146 spectra were then randomly split into training (70%) and test (30%) sets, and a supervised  
15  
16 147 classification model was constructed using support vector machines (SVM) to systemically predict  
17  
18 148 the time-after-harvest regardless the storage condition. The training set was used to build the SVM  
19  
20 149 training model, and the test set to evaluate its predictive ability. The SVM model was optimised by  
21  
22 150 cross-validation venetian blinds with 10 data splits.

23  
24  
25  
26 151 Metrics such as accuracy, sensitivity and specificity were calculated for the test set. For more than  
27  
28 152 two-classes, these metrics are calculated individually per class; herein, the average for all classes is  
29  
30 153 reported. The accuracy (AC) represents the total number of samples correctly classified considering  
31  
32 154 true and false negatives; the sensitivity (SENS) represents the proportion of positives that are  
33  
34 155 correctly classified; and the specificity (SPEC) represents the proportion of negatives that are  
35  
36 156 correctly identified <sup>28</sup>. These metrics are calculated as follow:

37  
38  
39  
40 157 
$$AC(\%) = \left( \frac{TP + TN}{TP + FP + TN + FN} \right) \times 100 \quad (1)$$

41  
42  
43  
44 158 
$$SENS(\%) = \left( \frac{TP}{TP + FN} \right) \times 100 \quad (2)$$

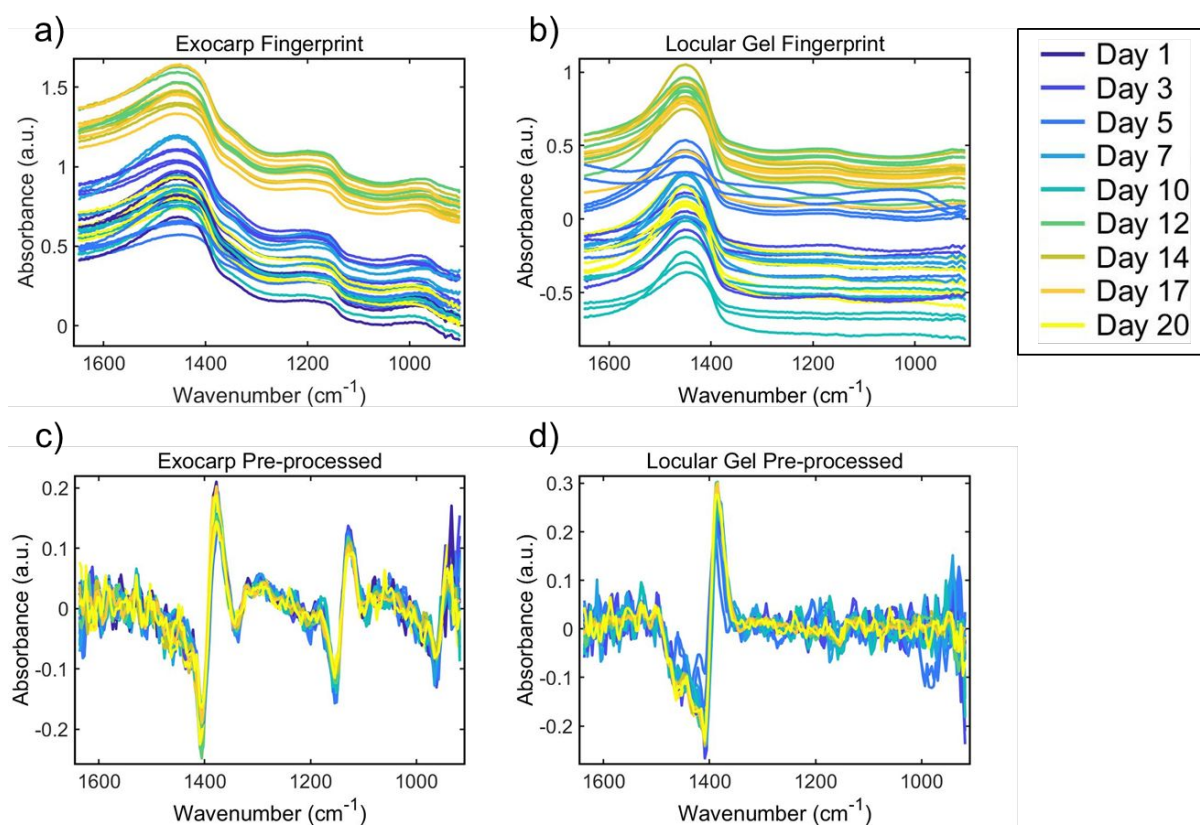
45  
46  
47  
48 159 
$$SPEC(\%) = \left( \frac{TN}{TN + FP} \right) \times 100 \quad (3)$$

49  
50  
51 160 where TP stands for true positives, TN for true negatives, FP for false positives, and FN for false  
52  
53 161 negatives.



162 **Results**

163 Figure 1 shows the fingerprint spectra for a) exocarp and b) locular gel absorbances, and pre-  
164 processed spectra for c) exocarp and d) locular gel. The pre-processed spectra for the exocarp  
165 measurements in Figure 1c contain bands around 1000 nm (C-H stretching 3<sup>rd</sup> overtone), 1200 nm  
166 (C-H stretching 2<sup>nd</sup> overtone in fibre parameters such as cellulose and lignin), 1360 nm (small arm, R-  
167 O-H stretching 1<sup>st</sup> overtone in alcohol), 1450 nm (O-H stretching 1<sup>st</sup> overtone in water), and spectral  
168 differences at around 1700 nm (C-H stretching 1<sup>st</sup> overtone in glucose/lignin)<sup>29,30</sup>. The locular gel  
169 fingerprint spectra, Figure 1d, have their bands compressed by the strong water peak at 1450 nm<sup>23</sup>.



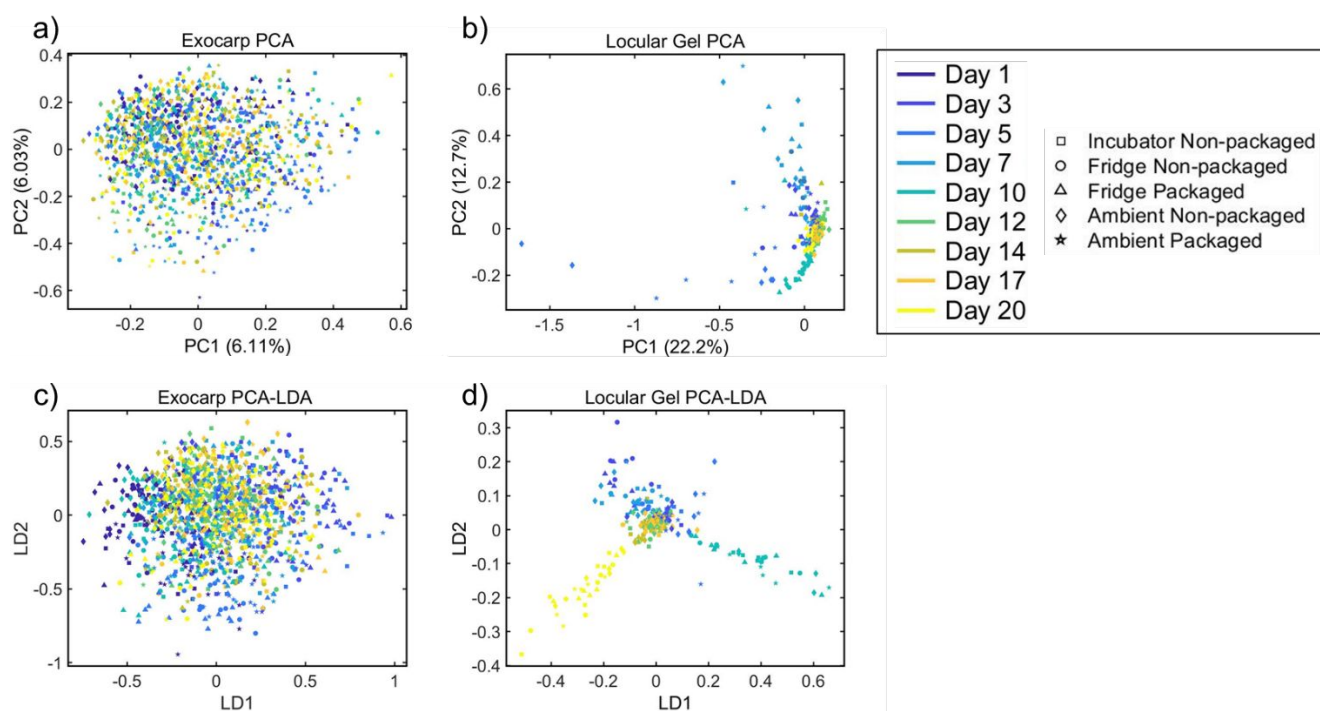
170

171 **Figure 1:** Fingerprint spectra for a) exocarp and b) locular gel absorbances, and pre-processed  
172 spectra for c) exocarp and d) locular gel absorbances. Each line is a class mean of spectra from a  
173 specific day and treatment condition. The colours represent the day the spectra were taken.

174 The pre-processed spectral data initially underwent unsupervised exploratory analysis by PCA,  
175 where overall segregation trends were observed in the data. Figure 2 shows PCA scatter plots for a)

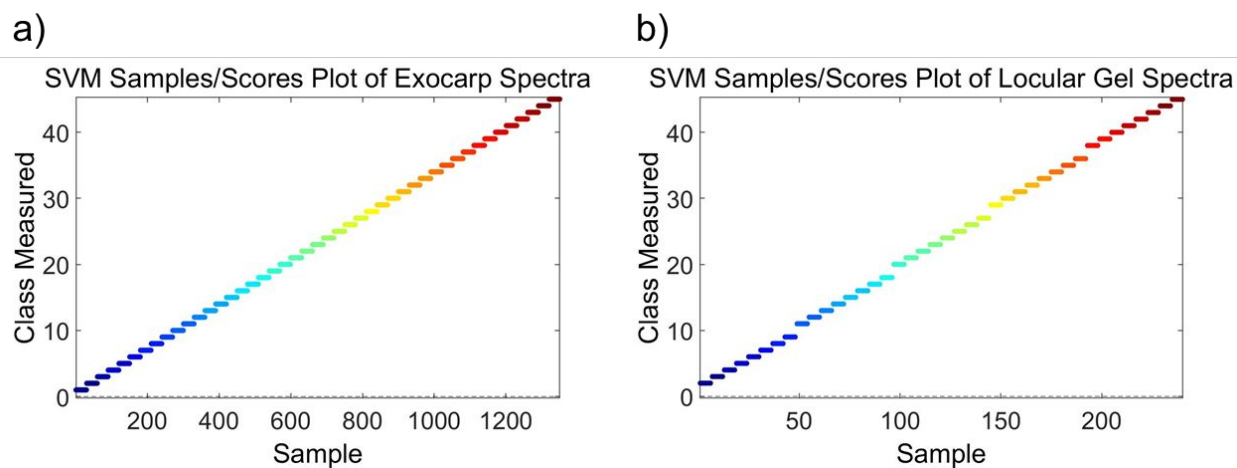
176 exocarp and b) locular gel spectra. The colours represent the time in days and the marker shapes represent the treatment and numbers inside parenthesis represent the explained variance for each PC. PCA did not provide separation between the samples (Figures 2a and 2b). For PCA scores of each storage condition individually, see supporting information Figure S2 and S3, for exocarp and locular gel measurements respectively.

181 Following the initial PCA, a supervised method, linear discriminant analysis, was also applied. PCA-LDA scatter plots are shown in Figure 2c) for exocarp and 2d) for locular gel spectra. Although there is some evidence of a time evolution trend for the exocarp measurements (see Figure 2c), the samples cannot be differentiated into clear clusters. Some clustering was however achieved for locular gel samples. Figure 2d shows that days 1, 10 and 20 are separated out best in the locular gel PCA-LDA scatter plot along the axes LD1 and LD2.



187  
188 **Figure 2:** PCA scatter plots for a) exocarp and b) locular gel spectra. Numbers inside parenthesis  
189 represent the explained variance for each PC. PCA-LDA scatter plots for c) exocarp and d) locular gel  
190 spectra. The colours represent the time in days and the marker shapes represent the treatment.

1  
2  
3 191 Figure 3 shows the SVM model prediction results for each time and treatment sample type grouped  
4  
5 192 separately. Overall, exocarp achieved higher accuracy sensitivity and specificity than locular gel  
6  
7 193 spectra. For all forty-five exocarp and forty-one local gel categories, the SVM model achieved poor  
8  
9 194 results at 51% and 54% accuracy, respectively. The SVM model to predict the time-after-harvest  
10  
11 195 regardless the storage condition for the exocarp spectra achieved 92% accuracy, 86% sensitivity and  
12  
13 196 98% specificity in test sets. For locular gel spectra, the SVM model achieved 84% accuracy, 74%  
14  
15 197 sensitivity and 95% specificity in test sets. Sensitivity of the exocarp spectra was - high, with many  
16  
17 198 samples correctly classified. For locular gel samples, packaged fridge-stored samples at days 3 and  
18  
19 199 20, and ambient stored non-packaged samples at day 3 achieved higher true positive rates than the  
20  
21 200 other groups. To give an overall picture of the trend over time, despite treatment, spectra from  
22  
23 201 different storage conditions were also grouped together. The SVM results for this grouping can be  
24  
25 202 viewed in the supporting information: see Figure S4 for optimisation parameters, Table S1 for  
26  
27 203 classification metrics, and Table S2 for test set confusion matrices.



204  
205 **Figure 3:** SVM predicted results for the test set for a) exocarp and b) locular gel. For all forty-five  
206 exocarp categories, the SVM model achieved 92% accuracy, 86% sensitivity and 98% specificity in  
207 test sets. For all forty-one locular gel categories, the SVM model achieved 84% accuracy, 74%  
208 sensitivity and 95% specificity in test sets. Each predicted category is shown in a different colour.

209 PCA loadings were used to determine the key wavenumbers associated with biochemical changes  
 210 over time for each treatment, in both exocarp and locular gel spectra (see Tables 1 and 2,  
 211 respectively). These spectral changes related to their biological origin by seeking molecular  
 212 assignments from existing literature. Spectra from each treatment group were considered separately  
 213 so that the changes highlighted in Tables 1 and 2 relate only to time.

214 **Table 1:** Loadings showing the key wavenumber changes for locular gel spectra associated with  
 215 biochemical changes over time for each treatment

Sample Type	Treatment Type	Loadings	Assignment	Reference
Locular Gel	INP	1636	CHO stretching of carbonyl group, typical saccharide absorption	30
		1456	CH <sub>3</sub> bending vibration (lipids and proteins)	30
		1170	C-O bands from glycomaterials and proteins	30
		1095	Stretching PO <sub>2</sub> symmetric	30
		981	Phosphodiester region	30
		918	Polysaccharides	31
Locular Gel	FP	1605	$\nu_{as}(\text{COO}^-)$ (polysaccharides, pectin)	30
		1571	C=N adenine	30
		1381	Amide II	32
		1075	Symmetric phosphate stretching modes or $\nu(\text{PO}_2^-)$ sym.	30
		1012	Starch	32
		945	D-(-)-Arabinose	33
Locular Gel	ANP	1636	C=O stretching of carbonyl group, typical saccharide absorption	30
		1596	Methylated nucleotides	30
		1443	$\delta(\text{CH})$ (polysaccharides, pectin), $\delta(\text{CH}_2)$ , lipids, fatty acids	30
		1170	C-O bands from glycomaterials and proteins	30
		993	Arabinoxylans	34
		930	Polysaccharides	31
Locular Gel	AP	1636	C=O stretching of carbonyl group, typical saccharide absorption	30
		1493	Protein	35
		1422	Protein and lipids	35
		1124	Polysaccharides	36
		977	C-O-C stretching at the $\beta$ -(1 $\rightarrow$ 4)-glycosidic linkages of amorphous cellulose	37
		918	Polysaccharides	31

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

<i>Locular Gel</i>	FNP	1609	Adenine vibration in DNA	30
		1574	C=N adenine	30
		1500	In-plane CH bending vibration from the phenyl rings, or Amide II (an N-H bending vibration coupled to C-N stretching)	30
		1236	Amide III and asymmetric phosphodiester stretching mode, $\nu_{as}(\text{PO}_2^-)$ mainly from the nucleic acids	30
		993	Arabinoxylans	34
		926	Polysaccharides	31

216

217

218 **Table 2:** Loadings showing the key wavenumber changes for exocarp spectra associated with  
 219 biochemical changes over time for each treatment

Sample Type	Treatment Type	Loadings	Assignment	Reference
Exocarp	INP	1636	C=O stretching of carbonyl group, typical saccharide absorption	30
		1590	C=O stretching of aliphatic and aromatic carbonyl groups	30
		1542	Amide II	30
		1483	Protein and lipids	35
		962	Polysaccharides	31
		918	Polysaccharides	31
Exocarp	FP	1627	Phenolic compounds	38
		1584	Amide II	30
		1526	C=N guanine	30
		1456	CH <sub>3</sub> bending vibration (lipids and proteins)	30
		981	Phosphodiester region	30
		918	Polysaccharides	31
Exocarp	ANP	1636	C=O stretching of carbonyl group, typical saccharide absorption	30
		1574	C=N adenine	30
		1532	Stretching C=N, C=C	30
		1479	Amino acid; $\nu[\text{COO}^-]$	32
		1381	Amide II	32
		918	Polysaccharides	31
Exocarp	AP	1627	Phenolic compounds	38
		1581	Ring C-C stretch of phenyl	30
		1509	Hemicellulose, C-C and C=C	39
		1473	Glycerolipids, wax hydrocarbons, $\delta(\text{CH}_2)$ scissoring	14
		1409	Succinic acid	40
		933	Z type DNA	30
Exocarp	FNP	1630	Amide I	30
		1593	Protein	36
		1509	Hemicellulose, C-C and C=C	39
		1436	Phenolic compounds	38
		1385	Structural polysaccharides, cellulose, C-H bending	41
		926	Polysaccharides	31

220

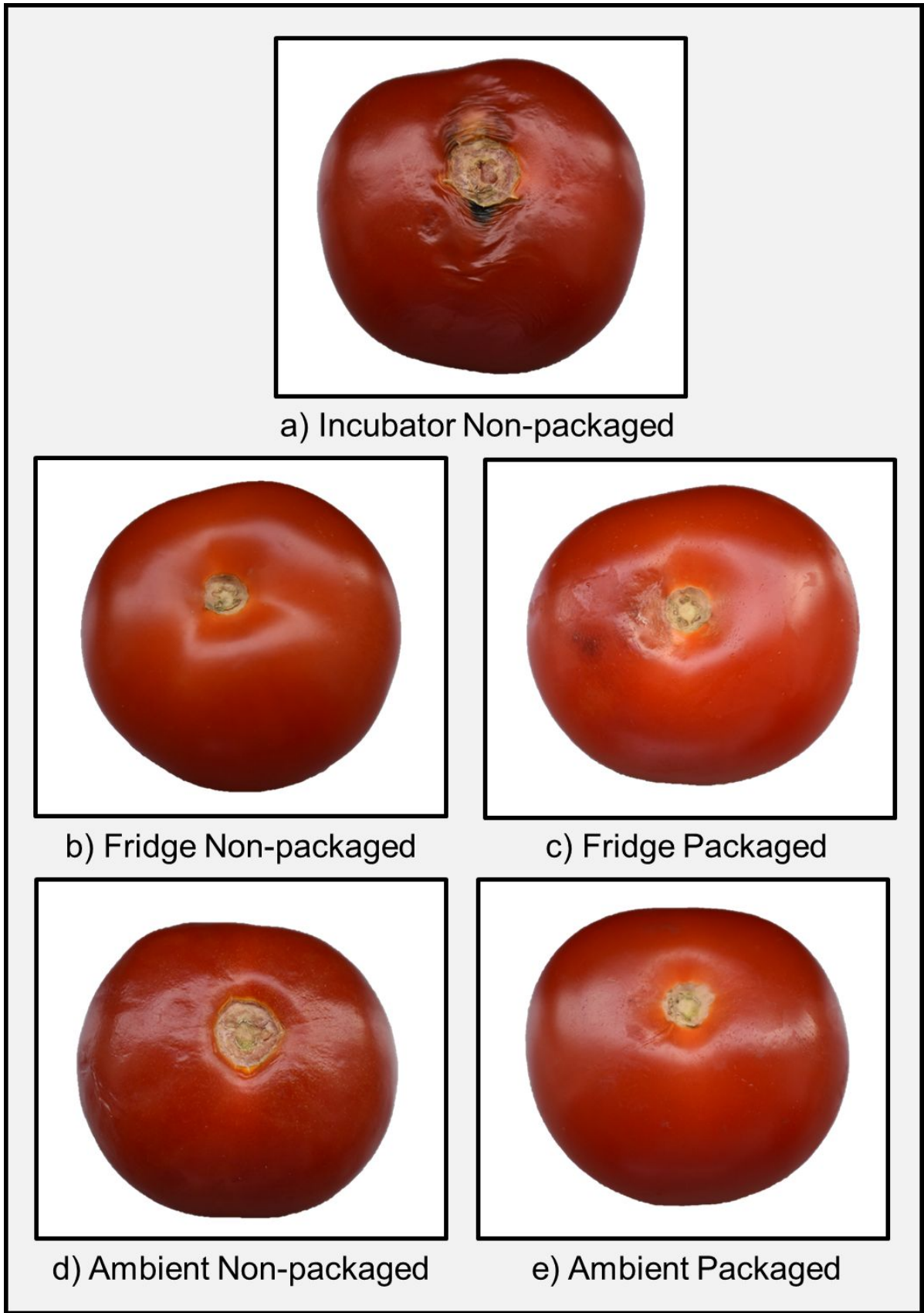
221 To compare the effects of different storage treatments on spectral absorbances, data from the  
 222 middle of the time course prior to any visible change in tomato quality were selected. Spectra from

day 10 were used to derive loadings showing which wavenumbers were associated with differing treatment conditions; these were subsequently connected with their corresponding biomarkers. To view loading graphs from which the information in Tables 1-3 were derived, see the supporting information Figure S5 for locular gel, Figure S6 for exocarp and Figure S7 for day ten.

**Table 3:** Loadings showing the key wavenumber differences between treatments at day 10 for locular gel and exocarp spectra

Sample Type	Day	Wavenumber	Molecular Assignment	Reference
<i>Exocarp</i>	Day 10	1633	Cutan, carboxylate functional groups ( $\nu_a(\text{COO}^-)$ )	38
		1590	C=O stretching of aliphatic and aromatic carbonyl groups	30
		1526	C=N guanine	30
		1479	Amino acid; $\nu[\text{COO}^-]$	32
		1436	Phenolic compounds	38
		1385	Structural polysaccharides, cellulose, C-H bending	41
<i>Locular Gel</i>	Day 10	1602	Pectin, phenolic compounds; C-C aromatic stretching; C-O-O <sup>-</sup> asymmetric stretching	42
		1486	Proteins	43
		1440	Phenolic compounds; $\nu(\text{C-C})$ aromatic (conjugated with C=C)	38
		1064	C-O stretching, C-C stretching (mannose-containing hemicellulose)	33
		993	Arabinoxylans	34
		933	Z type DNA	30

The different treatment conditions resulted in tomatoes of different qualities after 20 days (Figure 4). Visually, fridge non-packaged and ambient packaged treatments best preserved the quality of the tomatoes, whereas tomatoes from the incubator non-packaged and ambient non-packaged treatments showed signs of infection, and fridge packaged tomatoes showed signs of a possible chilling stress.



**Figure 4:** Images illustrating relative tomato quality between treatments on day 20



1  
2  
3 238 **Discussion**  
4

5 239 SVM achieved outstanding classification accuracy across exocarp (92%) and locular gel (84%)  
6  
7 240 samples to estimate the time-after-harvest despite the storage condition. The accuracy to predict all  
8  
9 241 forty-five and forty-one categories of samples (including differences between storage conditions)  
10  
11 242 could be improved if the model was built to differentiate between fewer categories, or a larger  
12  
13 243 dataset was used. When more than ten classes are involved in the classification model its  
14  
15 244 performance tends to reduce substantially since the dataset size is relatively reduced by a fraction of  
16  
17 245 the number of classes when building the classifier on a one-against-the-others categories basis [ref.]  
18  
19 246 Therefore, a very large dataset would be required to perform this type of classification. The models  
20  
21 247 were better at identifying where not to place a spectrum but often failed to find the correct category  
22  
23 248 to sort it into, with high specificities (60% and 91%) but low sensitivities (3% and 20%). When spectra  
24  
25 249 from all storage condition treatments were grouped together, classification by time achieved  
26  
27 250 excellent results with SVM; 92% average accuracy, 86% average sensitivity, and 98% average  
28  
29 251 specificity in the test set for exocarp measurements, and 84% average accuracy, 74% average  
30  
31 252 sensitivity, and 95% average specificity for the test set of locular gel samples, see supporting  
32  
33 253 information Table S2.  
34  
35  
36  
37  
38

39 254 PCA loadings identified key wavenumbers which allow discrimination between spectral absorbances  
40  
41 255 from tomatoes of different storage conditions and maturities, see Tables 1, 2 and 3. These were  
42  
43 256 connected to biomarkers using existing literature databases, allowing insight into the biochemical  
44  
45 257 changes taking place. Many of the identified changes related to the progression of starch  
46  
47 258 degradation, where the starch stored during development and is converted to soluble sugars, a key  
48  
49 259 process during post-harvest fruit ripening <sup>11</sup>. Wavenumber 1012 cm<sup>-1</sup> in locular gel fridge-stored  
50  
51 260 spectra identified starch specifically as an important indicator of tomato age <sup>32</sup>. All treatment  
52  
53 261 conditions for both locular gel and exocarp spectra identified that peaks for polysaccharides <sup>31</sup> and/  
54  
55 262 or saccharides <sup>30</sup> were used for the differentiation between tomatoes of different ages. This  
56  
57  
58  
59  
60

1  
2  
3 263 indicates that all fruits underwent some level of post-harvest ripening, irrespective of storage  
4  
5 264 conditions.

6  
7  
8 265 Other key peaks identified across different conditions were those relating to the structural and  
9  
10 266 compositional development of the cuticle and cell wall. Compositional changes in key compounds  
11  
12 267 such as pectin, cellulose, and other polysaccharides, as well as changes in cell wall thickness, are part  
13  
14 268 of the ripening process <sup>10</sup>. Tomato maturity was indicated in the spectra of ambient packaged  
15  
16 269 exocarp samples at peak 1473 cm<sup>-1</sup>, which relates to the  $\delta(\text{CH}_2)$  scissoring of cuticle glycerolipids and  
17  
18 270 wax hydrocarbons <sup>14</sup>. Pectin was identified as a measure of tomato maturity in ambient non-  
19  
20 271 packaged and fridge packaged samples from locular gel spectra <sup>30</sup>, and also as a differentiator  
21  
22 272 between storage conditions <sup>42</sup>. Cellulose was an indicator of tomato maturity in the locular gel of  
23  
24 273 ambient packaged fruit <sup>37</sup>. However cellulose and hemicellulose was more commonly associated  
25  
26 274 with exocarp spectra where these compounds differentiated between maturity in ambient packaged  
27  
28 275 and fridge non-packaged samples <sup>39</sup>. Cellulose was a key differentiator between treatment  
29  
30 276 conditions in exocarp spectra <sup>41</sup>.

31  
32  
33  
34  
35  
36 277 Dissolution of the cell wall is key to another ripening process, softening. Arabinogalactan proteins,  
37  
38 278 which act as a cross-linker for pectin and arabinoxylan, are thought to be important in the alteration  
39  
40 279 of cell wall assembly <sup>44</sup>. In this study, locular gel arabinoxylan levels were a key differentiator  
41  
42 280 between treatment conditions and allowed maturity determination in fridge and ambient non-  
43  
44 281 packaged fruit <sup>34</sup>. These results suggest that the packaging is altering the environment experienced  
45  
46 282 by the tomatoes, and having an impact on arabinoxylan, and consequently fruit-softness. The  
47  
48 283 texture changes can be seen in Figure 4, where the ambient non-packaged fruit appears softer than  
49  
50 284 its packaged equivalent. Expression of arabinogalactan proteins in tomato fruit has also been linked  
51  
52 285 to possible involvement in stress adaptations <sup>44</sup>, which may be why the packaged tomatoes suffered  
53  
54 286 from a chilling stress in the fridge whilst their unpacked equivalents appeared to be in better  
55  
56 287 condition, see Figure 4. Fleshy fruit such as tomatoes are particularly vulnerable to chilling injury, is a  
57  
58  
59  
60

1  
2  
3 288 type of oxidative stress that occurs during storage below 10 °C <sup>45</sup>. Despite this chilling, fridge non-  
4  
5 289 packaged tomatoes were amongst the best-preserved tomatoes alongside ambient packaged  
6  
7 290 tomatoes, see Figure 4. The least well-preserved tomatoes were those not protected with packaging,  
8  
9 291 stored under ambient room temperature or in the incubator, which showed signs of infection. The  
10  
11 292 lack of packaging likely exposed the tomatoes to any pathogens present, and the high temperatures  
12  
13 293 in the incubator conditions promoted growth. Bacteria and fungi grow best at temperatures  
14  
15 294 between 25-30°C <sup>46,47</sup>, and the incubated tomatoes were stored at 25°C.

16  
17  
18  
19 295 This study has identified biomarkers, including those indicating tomato ripening, which are modified  
20  
21 296 by storage and packaging conditions and that have the potential to be used to target reductions in  
22  
23 297 the unnecessary disposal of tomatoes through spoilage in the supply chain. The rapid and non-  
24  
25 298 destructive nature of exocarp scanning and the portability of the spectrometer, capable of  
26  
27 299 connection to a mobile phone, renders this method particularly suitable for use within the food  
28  
29 300 industry. With further development, this user-friendly and non-destructive technology displays  
30  
31 301 potential for a wider range of applications within and beyond the food industry.

### 32 33 34 35 36 302 **Acknowledgements**

37  
38 303 We would like to acknowledge John Lane, who supplied the tomatoes and Lancaster Girls' Grammar  
39  
40 304 School for their support for the project.

### 41 42 43 44 305 **Author Contributions**

45  
46 306 NEME, FLM and CLMM conceived the study; NEME, SG, RSB, FLM and CLMM jointly designed the  
47  
48 307 experiment; NEME, SG, RSB conducted the experiment and collected the data; CLMM and CAH  
49  
50 308 undertook the chemometric analyses; NEME, CAH, SG, RSB, CLMM jointly drafted the manuscript,  
51  
52 309 MRM advised on the plant biology and all authors approved the final version; CR provided support  
53  
54 310 and supervision for the conduct of the experiment; FLM was principal investigator.

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

311 **Declaration of Interest**

312 FLM is a major shareholder in Biocel UK Ltd, a company with interests in developing tools described  
313 herein for commercial gain.

314

315 **References**

- 316 1. Food and Agriculture Organization of the United Nations. High-level expert forum—How to  
317 feed the world in 2050. (2009).
- 318 2. Nino, F. S. Sustainable Development Goals—United Nations. *United Nations Sustain. Dev.*  
319 (2015).
- 320 3. Parfitt, J., Barthel, M. & MacNaughton, S. Food waste within food supply chains:  
321 Quantification and potential for change to 2050. *Philos. Trans. R. Soc. B Biol. Sci.* **365**, 3065–  
322 3081 (2010).
- 323 4. Story, E. N., Kopec, R. E., Schwartz, S. J. & Harris, G. K. An Update on the Health Effects of  
324 Tomato Lycopene. <http://dx.doi.org/10.1146/annurev.food.102308.124120> **1**, 189–210  
325 (2010).
- 326 5. Wu, X., Yu, L. & Pehrsson, P. R. Are Processed Tomato Products as Nutritious as Fresh  
327 Tomatoes? Scoping Review on the Effects of Industrial Processing on Nutrients and Bioactive  
328 Compounds in Tomatoes. *Adv. Nutr.* (2021) doi:10.1093/ADVANCES/NMAB109.
- 329 6. Boz, Z. & Sand, C. K. A systematic analysis of the overall nutritional contribution of food loss  
330 and waste in tomatoes, spinach, and kidney beans as a function of processing. *J. Food Process*  
331 *Eng.* **43**, e13509 (2020).
- 332 7. Xue, L. *et al.* Mapping the EU tomato supply chain from farm to fork for greenhouse gas  
333 emission mitigation strategies. *J. Ind. Ecol.* **25**, 377–389 (2021).
- 334 8. Ntinas, G. K., Neumair, M., Tsadilas, C. D. & Meyer, J. Carbon footprint and cumulative energy  
335 demand of greenhouse and open-field tomato cultivation systems under Southern and  
336 Central European climatic conditions. *J. Clean. Prod.* **142**, 3617–3626 (2017).
- 337 9. Qvested, T., Ingle, R. & Parry, A. Household food and drink waste in the United Kingdom  
338 2012, Waste & Resources Action Programme (WRAP). *Banbury, UK* (2013).
- 339 10. Segado, P., Domínguez, E. & Heredia, A. Ultrastructure of the Epidermal Cell Wall and Cuticle  
340 of Tomato Fruit (*Solanum lycopersicum* L.) during Development. *Plant Physiol.* **170**, 935

- 1  
2  
3 341 (2016).  
4  
5 342 11. Fan, Z.-Q. *et al.* A banana R2R3-MYB transcription factor MaMYB3 is involved in fruit ripening  
6  
7 343 through modulation of starch degradation by repressing starch degradation-related genes  
8  
9 344 and MabHLH6. *Plant J.* **96**, 1191–1205 (2018).  
10  
11 345 12. Skolik, P., McAinsh, M. R. & Martin, F. L. ATR-FTIR spectroscopy non-destructively detects  
12  
13 346 damage-induced sour rot infection in whole tomato fruit. *Planta* **249**, 925–939 (2019).  
14  
15 347 13. Clément, A., Dorais, M. & Vernon, M. Multivariate Approach to the Measurement of Tomato  
16  
17 348 Maturity and Gustatory Attributes and Their Rapid Assessment by Vis–NIR Spectroscopy. *J.*  
18  
19 349 *Agric. Food Chem.* **56**, 1538–1544 (2008).  
20  
21 350 14. Skolik, P., Morais, C. L. M., Martin, F. L. & McAinsh, M. R. Determination of developmental  
22  
23 351 and ripening stages of whole tomato fruit using portable infrared spectroscopy and  
24  
25 352 Chemometrics. *BMC Plant Biol.* **19**, 236 (2019).  
26  
27 353 15. Tiwari, G., Slaughter, D. C. & Cantwell, M. Nondestructive maturity determination in green  
28  
29 354 tomatoes using a handheld visible and near infrared instrument. *Postharvest Biol. Technol.*  
30  
31 355 **86**, 221–229 (2013).  
32  
33 356 16. Huang, Y., Lu, R. & Chen, K. Prediction of firmness parameters of tomatoes by portable visible  
34  
35 357 and near-infrared spectroscopy. *J. Food Eng.* **222**, 185–198 (2018).  
36  
37 358 17. Sirisomboon, P., Tanaka, M., Kojima, T. & Williams, P. Nondestructive estimation of maturity  
38  
39 359 and textural properties on tomato ‘Momotaro’ by near infrared spectroscopy. *J. Food Eng.*  
40  
41 360 **112**, 218–226 (2012).  
42  
43 361 18. Saad, A., Gawad Saad, A., Jaiswal, P. & Narayan Jha, S. Non-destructive quality evaluation of  
44  
45 362 intact tomato using VIS-NIR spectroscopy Related papers Non-destructive quality evaluation  
46  
47 363 of intact tomato using VIS-NIR spectroscopy. *Int. J. Adv. Res.* **2**, 632–639 (2014).  
48  
49 364 19. Nolasco Perez, I. *et al.* Classification of Chicken Parts Using a Portable Near-Infrared (NIR)  
50  
51 365 Spectrophotometer and Machine Learning. *Appl. Spectrosc.* **72**, 1774–1780 (2018).  
52  
53 366 20. Kaufmann, K. *et al.* Portable NIR Spectrometer for Prediction of Palm Oil Acidity. *J. Food Sci.*  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 367 **84**, 406–411 (2019).  
4  
5 368 21. Oliveira, M. M., Cruz-Tirado, J. P., Roque, J. V, Teófilo, R. F. & Barbin, D. F. Portable near-  
6  
7 369 infrared spectroscopy for rapid authentication of adulterated paprika powder. *J. Food*  
8  
9 370 *Compos. Anal.* **87**, 103403 (2020).  
10  
11 371 22. Beć, K. B., Grabska, J. & Huck, C. W. Principles and Applications of Miniaturized Near-Infrared  
12  
13 372 (NIR) Spectrometers. *Chem. – A Eur. J.* **27**, 1514–1532 (2021).  
14  
15 373 23. Morais, C. L. M., Lima, K. M. G., Singh, M. & Martin, F. L. Tutorial: multivariate classification  
16  
17 374 for vibrational spectroscopy in biological samples. *Nature Protocols* vol. 15 2143–2162 (2020).  
18  
19 375 24. Atta-Aly, M. A., Brecht, J. K. & Huber, D. J. Ripening of tomato fruit locule gel tissue in  
20  
21 376 response to ethylene. *Postharvest Biol. Technol.* **19**, 239–244 (2000).  
22  
23 377 25. Bro, R. & Smilde, A. K. Principal component analysis. *Anal. methods* **6**, 2812–2831 (2014).  
24  
25 378 26. Morais, C. L. M. & Lima, K. M. G. Principal Component Analysis with Linear and Quadratic  
26  
27 379 Discriminant Analysis for Identification of Cancer Samples Based on Mass Spectrometry. *Artic.*  
28  
29 380 *J. Braz. Chem. Soc* **29**, 472–481 (2018).  
30  
31 381 27. Cortes, C. & Vapnik, V. Support-vector networks. *Mach. Learn.* **20**, 273–297 (1995).  
32  
33 382 28. Morais, C. L. M. & Lima, K. M. G. Comparing unfolded and two-dimensional discriminant  
34  
35 383 analysis and support vector machines for classification of EEM data. *Chemom. Intell. Lab.*  
36  
37 384 *Syst.* **170**, 1–12 (2017).  
38  
39 385 29. Jerome, J. & Workman, J. Interpretive Spectroscopy for Near Infrared.  
40  
41 386 <http://dx.doi.org/10.1080/05704929608000571> **31**, 251–320 (2006).  
42  
43 387 30. Talari, A. C. S., Martinez, M. A. G., Movasaghi, Z., Rehman, S. & Rehman, I. U. Advances in  
44  
45 388 Fourier transform infrared (FTIR) spectroscopy of biological tissues. *Appl. Spectrosc. Rev.* **52**,  
46  
47 389 456–506 (2017).  
48  
49 390 31. Di Giambattista, L. *et al.* New marker of tumor cell death revealed by ATR-FTIR spectroscopy.  
50  
51 391 *Anal. Bioanal. Chem.* **399**, 2771–2778 (2011).  
52  
53 392 32. Luo, Y. *et al.* Diagnostic Segregation of Human Breast Tumours Using Fourier-transform  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 393 Infrared Spectroscopy Coupled with Multivariate Analysis: Classifying Cancer Subtypes.  
4  
5 394 *Spectrochim. Acta Part A Mol. Biomol. Spectrosc.* **255**, 119694 (2021).  
6  
7  
8 395 33. Liu, X., Renard, C. M. G. C., Bureau, S. & Le Bourvellec, C. Revisiting the contribution of ATR-  
9  
10 396 FTIR spectroscopy to characterize plant cell wall polysaccharides. *Carbohydr. Polym.* **262**,  
11  
12 397 117935 (2021).  
13  
14 398 34. Demir, P., Onde, S. & Severcan, F. Phylogeny of cultivated and wild wheat species using ATR-  
15  
16 399 FTIR spectroscopy. *Spectrochim. Acta - Part A Mol. Biomol. Spectrosc.* **135**, 757–763 (2015).  
17  
18  
19 400 35. Meinen, C. & Rauber, R. Root discrimination of closely related crop and weed species using FT  
20  
21 401 MIR-ATR spectroscopy. *Front. Plant Sci.* **6**, 765 (2015).  
22  
23 402 36. Butler, H. J., McAinsh, M. R., Adams, S. & Martin, F. L. Application of vibrational spectroscopy  
24  
25 403 techniques to non-destructively monitor plant health and development. *Anal. Methods* **7**,  
26  
27 404 4059–4070 (2015).  
28  
29  
30 405 37. Bekiaris, G. *et al.* Rapid estimation of sugar release from winter wheat straw during  
31  
32 406 bioethanol production using FTIR-photoacoustic spectroscopy. *Biotechnol. Biofuels* **2015** **8**,  
33  
34 407 1–12 (2015).  
35  
36  
37 408 38. Heredia-Guerrero, J. A. *et al.* Infrared and Raman spectroscopic features of plant cuticles: a  
38  
39 409 review. *Front. Plant Sci.* **0**, 305 (2014).  
40  
41 410 39. Zhou, C., Jiang, W., Via, B. K., Fasina, O. & Han, G. Prediction of mixed hardwood lignin and  
42  
43 411 carbohydrate content using ATR-FTIR and FT-NIR. *Carbohydr. Polym.* **121**, 336–341 (2015).  
44  
45  
46 412 40. Kang, S., Amarasiriwardena, D. & Xing, B. Effect of dehydration on dicarboxylic acid  
47  
48 413 coordination at goethite/water interface. *Colloids Surfaces A Physicochem. Eng. Asp.* **318**,  
49  
50 414 275–284 (2008).  
51  
52  
53 415 41. Rana, R. *et al.* Leaf Attenuated Total Reflection Fourier Transform Infrared (ATR-FTIR)  
54  
55 416 biochemical profile of grassland plant species related to land-use intensity. *Ecol. Indic.* **84**,  
56  
57 417 803–810 (2018).  
58  
59 418 42. Moskal, P., Weselucha-Birczyńska, A., Łabanowska, M. & Filek, M. Adaxial and abaxial pattern  
60



- 1  
2  
3 419 of *Urtica dioica* leaves analyzed by 2DCOS ATR-FTIR as a function of their growth time and  
4  
5 420 impact of environmental pollution. *Vib. Spectrosc.* **104**, 102948 (2019).  
6  
7  
8 421 43. Jin, N. *et al.* Spectrochemical analyses of growth phase-related bacterial responses to low  
9  
10 422 (environmentally-relevant) concentrations of tetracycline and nanoparticulate silver. *Analyst*  
11  
12 423 **143**, 768–776 (2018).  
13  
14 424 44. Leszczuk, A., Kalaitzis, P., Blazakis, K. N. & Zdunek, A. The role of arabinogalactan proteins  
15  
16 425 (AGPs) in fruit ripening—a review. *Hortic. Res.* **2020 71 7**, 1–12 (2020).  
17  
18  
19 426 45. Stevens, R. *et al.* Tomato fruit ascorbic acid content is linked with monodehydroascorbate  
20  
21 427 reductase activity and tolerance to chilling stress. *Plant. Cell Environ.* **31**, 1086–1096 (2008).  
22  
23 428 46. Jay, J. M., Loessner, M. J. & Golden, D. A. Low-temperature food preservation and  
24  
25 429 characteristics of psychrotrophic microorganisms. in *Modern food microbiology* (Springer  
26  
27 430 Science & Business Media, 2008).  
28  
29  
30 431 47. Adams, M. R. & Moss, M. O. The microbiology of food preservation. (2020).  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60