

Document downloaded from:

<http://hdl.handle.net/10251/156015>

This paper must be cited as:

Ibañez, G.; Valcárcel-Germes, M.; Cebolla Cornejo, J.; Rosello Ripolles, S. (2019). FT-MIR determination of taste-related compounds in tomato: a high throughput phenotyping analysis for selection programs. *Journal of the Science of Food and Agriculture*. 99(11):5140-5148. <https://doi.org/10.1002/jsfa.9760>



The final publication is available at

<https://doi.org/10.1002/jsfa.9760>

Copyright John Wiley & Sons

#### Additional Information

"This is the peer reviewed version of the following article: Ibáñez, Ginés, Mercedes Valcárcel, Jaime Cebolla-Cornejo, and Salvador Roselló. 2019. FT-MIR Determination of Taste-related Compounds in Tomato: A High Throughput Phenotyping Analysis for Selection Programs. *Journal of the Science of Food and Agriculture* 99 (11). Wiley: 5140 48. doi:10.1002/jsfa.9760, which has been published in final form at <https://doi.org/10.1002/jsfa.9760>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Self-Archiving."

1 **FT-MIR determination of taste-related compounds in tomato: a high throughput phenotyping**  
2 **analysis for selection programs**

3 **Running title: FT-MIR determination of taste-related compounds in tomato**

4

5 Ginés Ibáñez<sup>a</sup>, Mercedes Valcárcel<sup>a</sup>, Jaime Cebolla-Cornejo<sup>b\*</sup>, Salvador Roselló<sup>a</sup>

6 <sup>a</sup> Unidad Mixta de Investigación Mejora de la Calidad Agroalimentaria UJI-UPV. Department de  
7 Ciències Agràries i del Medi Natural, Universitat Jaume I, Avda. Sos Baynat s/n, 12071 Castelló  
8 de la Plana, Spain

9 <sup>b</sup> Unidad Mixta de Investigación Mejora de la Calidad Agroalimentaria UJI-UPV. COMAV.  
10 Universitat Politècnica de València, Cno. de Vera s/n, 46022 València, Spain

11 \*Corresponding author: jaicecor@btc.upv.es; Tel.: +34-963879423

12 ORCID codes: G. Ibáñez: 0000-0002-1787-8587; M. Valcárcel: 0000-0002-9347-1500; J. Cebolla-  
13 Cornejo: 0000-0002-2607-9920; S. Roselló: 0000-0002-7733-4178

14

15

16 **Abstract**

17 BACKGROUND: Tomato taste is defined by the accumulation of sugars and organic acids.  
18 Individual analyses of these compounds using HPLC or CZE are expensive, time-consuming and  
19 are not feasible for large number of samples, justifying the interest of spectroscopic methods  
20 such as Fourier transform mid-infrared (FT-MIR). This work analysed the performance of FT-MIR  
21 models to determine the accumulation of sugars and acids, considering the efficiency of models  
22 obtained with different ranges of variation.

23 RESULTS: FT-MIR spectra (five-bounce attenuated total reflectance, ATR) were used to obtain  
24 PLS models to predict sugar and acid contents in specific sample sets representing different  
25 varietal types. A general model was also developed, obtaining R<sup>2</sup> values for prediction higher  
26 than 0.84 for main components (SSC, fructose, glucose, and citric acid). Root mean squared error  
27 of prediction RMSEP for these components were lower than 15% of the mean contents and  
28 lower than 6% of the highest contents. Even more, the model sensitivity and specificity for those  
29 variables with a 10% selection pressure was 100%. That means that all samples with the 10%  
30 highest content were correctly identified. The model was applied to an external assay and it  
31 exhibited, for main components, high sensitivities (>70%) and specificities (>96%). RMSEP values  
32 for main compounds were lower than 21% and 13% of the mean and maximum content  
33 respectively.

34 CONCLUSION: The models obtained confirm the effectiveness of FT-MIR models to select  
35 samples with high contents of taste-related compounds, even when the calibration has not been  
36 performed within the same assay.

37

38 **Keywords:** *Solanum lycopersicum* L., sugars, organic acids, plant breeding, quality, FTIR

39

## 40 INTRODUCTION

41 Consumer complaints about tomato taste became commonplace during the nineties.<sup>1</sup> This  
42 discontent helped to consolidate emerging quality markets associated with tomato landraces  
43 with recognized organoleptic quality.<sup>2</sup> Several causes explain the loss of organoleptic quality in  
44 modern varieties, including the way tomatoes are produced (early harvesting, high nitrogen  
45 fertilization, greenhouse cultivation...), conserved (i.e. refrigeration) or selected in breeding  
46 programs.<sup>3-6</sup> This last cause had a strong influence. In part, it was due to negative collateral  
47 effects of genes controlling other interesting traits such as *uniform ripening* (*u*), which leads to  
48 reduced sugar contents in ripe fruits,<sup>7</sup> or *rin* and *nor*, which offer long shelf life but affect the  
49 production of flavor-related compounds.<sup>8,9</sup> Nevertheless, probably the main negative effect was  
50 due to the strong emphasis placed in the selection for high production and biotic stress  
51 resistance and the scant attention paid to fruit flavor.

52 Tomato flavor is determined by the accumulation of taste-related compounds and aroma  
53 volatiles. Among the first, taste depends on the accumulation of sugars, mainly fructose and  
54 glucose (with traces of sucrose), organic acids, mainly malic, citric and glutamic acids and the  
55 ratios between them. On the other hand, volatiles would not only affect aroma perception but  
56 also affect the way that tongue receptors perceive sweetness.<sup>6,10</sup>

57 In order to redirect breeding programs towards high-quality varieties, it would be necessary to  
58 determine the individual accumulation of sugars and acids in selection programs. Especially,  
59 when these new materials are targeted to cover segmented tomato markets, that demand  
60 quality profiles with subtle taste differences.<sup>11</sup> For this purpose, taste-related compounds such  
61 as sugars and acids can be determined by precise instrumental techniques such as liquid  
62 chromatography, HPLC,<sup>12</sup> or capillary zone electrophoresis, CZE<sup>13</sup>. However, these techniques  
63 have a high cost and require qualified personnel for their use. Even more, only a small number  
64 of samples can be processed per day considering the duration of analysis.

65 As an alternative, infrared spectroscopy can provide an indirect quantification of these  
66 compounds with several advantages. Its use requires minimum preparation of the sample, it is  
67 inexpensive and offers rapid analysis time.<sup>14</sup> Although Near-infrared (NIR) spectroscopy has  
68 been successfully used in measuring quality attributes of horticultural produces,<sup>15</sup> but  
69 sometimes the goodness of the models obtained retrain their application to general screening  
70 purposes<sup>16</sup> and studies in different contexts have evidenced a better performance of Fourier-  
71 transform mid-infrared (FT-MIR) spectroscopy in the quantification of sugars and acids.<sup>17</sup>

72 In tomato, FT-MIR models for the quantification of taste-related compounds are scarce. They  
73 have been obtained using a wide range of varieties<sup>18</sup> or a specific varietal type, such as  
74 processing tomato.<sup>19,20</sup> Although different FT-MIR methodologies have been tested, Wilkerson  
75 et al.<sup>19</sup> concluded that little difference was observed comparing the use of triple bounce  
76 attenuated total reflectance (ATR) and transmission. In their study, Wilkerson et al.<sup>19</sup> concluded  
77 that increasing the number and diversity of the samples would reduce the impact of irrelevant  
78 spectral-variations (noise) in the calibration model, thus resulting in partial least squared  
79 regression (PLSR) models with higher correlation coefficients.

80 In the present study, this premise is addressed: does a higher amplitude of samples and  
81 environments really increase the efficiency of FT-MIR models in the prediction of the  
82 concentration of taste-related compounds in tomato? At the same time, the feasibility of the  
83 use of five-bounce ATR FT-MIR robust models in selection procedures of fresh and processing  
84 tomato cultivars is studied. For that purpose, the prediction models obtained were tested with

85 an external assay, including 111 samples, representing the variation produced by changes in the  
86 environmental growing conditions.

87

## 88 **MATERIALS AND METHODS**

### 89 **Plant Material**

90 Three sets of samples were used to develop prediction models. Set number 1 included 108  
91 samples of processing tomato of eight varieties grown with different water and fertilization  
92 regimes in Extremadura (Spain). Set number 2 included 107 samples of medium sized tomatoes  
93 from 25 varieties including beef, rounded, plum and cluster tomatoes from commercial and  
94 traditional varieties. Set number 3 included 115 samples of 32 varieties of cherry and cocktail  
95 tomato. Samples from set 1 were obtained during the development of different agronomical  
96 studies.<sup>21,22</sup> Samples from sets 2 and 3 were purchased from local markets.

97 Each specific sample set and a general set grouping the total 330 tomato samples were used for  
98 the construction of models predicting sugar and acid contents from FT-MIR spectra. In all cases,  
99 fully ripe fruits were sampled.

100 A fourth sample set from a different external assay was also obtained. It contained 111 samples  
101 of processing tomato representing the same varieties and water regimes of sample set 1, but  
102 grown in Navarra (Spain) with different environmental conditions (lower radiation and  
103 temperature). Growing conditions were described with higher detail by Martí et al.<sup>22</sup>. These  
104 samples were not included in the calculation of the general model and were only used to test  
105 the robustness and true prediction capabilities of the model.

### 106 **Sample preparation**

107 The tomatoes were crushed, homogenized and stored at -80°C until analysis. Subsequently, the  
108 samples were thawed and centrifuged at 15680g for 10 minutes, following the conditions  
109 reported by Wilkerson et al. for FT-MIR analysis<sup>19</sup>. Three aliquots of the supernatant were  
110 obtained. One of them was used to determine soluble solids content (SSC) with a Pocket PAL- $\alpha$   
111 digital refractometer (Atago, Tokyo, Japan). The remaining two were used to obtain the FT-MIR  
112 spectra and sugar and acid contents via capillary zone electrophoresis (CZE).

### 113 **Infrared spectroscopy analysis**

114 The absorbance measurements of the FT-MIR spectrum were carried out using a portable Cary  
115 630 FT-MIR spectrometer (Agilent Technologies, Waldbronn, Germany), equipped with a DTGS  
116 (Deuterated triglycine sulfate) detector and a five-bounce ZnSe ATR. Microlab FT-MIR Software  
117 v. B.05.3. (Agilent Technologies, Waldbronn, Germany) was used to acquire the data, selecting  
118 a spectral resolution of 4 cm<sup>-1</sup> in a range of the average infrared spectrum of 4000-650 cm<sup>-1</sup>.  
119 After reviewing the spectra, and considering previous literature<sup>19</sup>, only the 1500-900 cm<sup>-1</sup>  
120 spectra were used for chemometrics.

121 Two independent spectral measurements (average of 64 consecutive scans) were performed in  
122 each sample. Measurements of the reference spectrum (background) were obtained between  
123 samples to correct uncontrolled variations in the spectral measurements due to variations in  
124 environmental conditions. Between the different measurements, the glass was carefully cleaned  
125 with distilled water and dried with cellulose tissues.

126 **CZE analysis**

127 The quantifications of the main reducing sugars (fructose and glucose) and organic acids (citric,  
128 malic and glutamic) were performed by capillary zone electrophoresis (CZE) with an Agilent 7100  
129 equipment (Agilent Technologies, Waldbronn, Germany) following the method described by  
130 Cebolla-Cornejo et al.<sup>13</sup>

131 **Data analysis**

132 Three specific prediction models for sample sets 1 to 3 and an additional general model with all  
133 these samples was calculated. Each model was calculated using 75% of the sample set  
134 (calibration group) to develop the calibration and cross-validation procedures. The remaining  
135 25% of the samples (validation group) were not included in the models and were used to obtain  
136 an accurate estimate of the error committed predicting the contents of new samples. Both the  
137 calibration and validation groups were randomly selected. In the case of the general model the  
138 calibration and validation groups were again randomly selected, independently of the selection  
139 made for each specific model.

140 FT-MIR spectra were pre-treated to eliminate multiplicative signal interferences with a  
141 Multiplicative Scatter Correction (MSC) and response variables were autoscaled using the mean  
142 and standard deviation. The predictive models were then obtained by least squares partial  
143 regression, PLS.<sup>23</sup> In order to choose the optimal number of latent variables, Venetian blinds  
144 cross-validation procedure was applied. In order to check the validity of the model root mean  
145 squared errors of calibration (RMSEC) and cross-validation (RMSECV) were calculated. RMSECV  
146 values were used as one of the selection criteria for the number of latent variables to be included  
147 in the model. In this sense, new latent variables were not included if they did not lead to a  
148 reduction of RMSECV higher than 2%. The second criterion used was to select the lowest  
149 possible number of latent variables.

150 Outliers in the FT-MIR spectra and response variables were identified and removed. The values  
151 of the Hotelling T2 statistics and the Q residues were considered for the former and the values  
152 of the normalized residuals (<-3 or >3) and leverage parameters were considered for the  
153 response variables. After calculating the PLS regression models, the FT-MIR spectra of the  
154 samples of the validation group were used to make predictions and root mean squared errors  
155 of prediction (RMSEP) were then calculated.

156 Correlation coefficients ( $R^2$ ) were calculated for the calibration ( $R^2_c$ ), cross-validation ( $R^2_{cv}$ ) and  
157 prediction data ( $R^2_p$ ). RMSEP values were also contextualized using the mean (%mean) and  
158 maximum (%maximum) values. Additionally, the predictive capacity of the models was assessed  
159 using the dimensionless parameter residual prediction deviation (RPD), which represents the  
160 ratio between the standard deviation of the validation and RMSEP, and the range error ratio  
161 (RER), which is the ratio between the range in the composition values of the validation samples  
162 and the RMSEP. RPD and RER enable a better comparison between models obtained with  
163 different samples, especially the former, as RER values are highly dependent on rare high  
164 contents. Usually, RPD values should be higher than 2 in order to represent useful models for  
165 classification or quantification.<sup>24,25</sup> On the other hand, Williams and Norris<sup>26</sup> suggested that RER  
166 > 10 highly useful models while they would have limited to good application for values between  
167 3 and 10. Additionally, the true prediction performance was assessed with external samples. In  
168 order to determine the effectiveness of the general model for screening, sensitivity (true  
169 positive rate) and specificity (true negative rate) values were calculated with the general model  
170 when a 10% or 20% selection pressure was applied.

171 The samples of the external model were used to determine if the obtained general model had a  
172 reliable performance not only with samples grown in the same environmental conditions  
173 (validation group), but also with those harvested in completely different conditions (external  
174 sample set).

175 Spectra pre-treatment, PLS regression models, detection of outliers, error parameters and  
176 goodness of fit for each model were performed in Matlab v 9.4 environment (Mathworks Inc,  
177 Natick, MA, USA ) using the PLS\_Toolbox v 8.2.1 module (Eigenvector Research Inc, Wenatchee,  
178 WA, USA). A heatmap of the correlations between statistical parameters of the samples and the  
179 performance of the PLS regression models was obtained with *Heatmapper*  
180 (<http://www.heatmapper.ca/pairwise/>).

181

## 182 **RESULTS**

183 The spectra obtained for the tomato samples were characterized by two high absorption areas  
184 corresponding to 3700-3000  $\text{cm}^{-1}$  and 1750-1500  $\text{cm}^{-1}$  (Fig. 1). In the first area, considerable  
185 differences were detected among samples, while in the second these differences were limited.  
186 Although absorbance peak in the area 1500-900  $\text{cm}^{-1}$  was not as high, important differences in  
187 this area were detected among samples, especially in the area 1150-950  $\text{cm}^{-1}$ .

### 188 ***Characterization of sample sets***

189 SSC values obtained were in agreement with the expected values considering the varietal types  
190 used. Fresh mid-sized tomatoes had the lowest value, followed by processing tomato and cherry  
191 and cocktail tomato (Table 1). Accordingly, fructose and glucose mean contents followed the  
192 same distribution. The contents of the acids (citric, malic and glutamic) had a different  
193 distribution. Cherry and cocktail maintained higher contents, but fresh mid-sized tomatoes had  
194 higher levels of acids compared to processing tomato. For all the models, the samples selected  
195 for the validation group, that were used to predict values using the calculated PLS models,  
196 represented a similar range of variation (Table 1).

197 The range of variation for each variable differed among sample sets (Table 1). In general, the  
198 range of variation for sugars was lower than for acids, with the lowest variation being present  
199 for SSC. Processing tomatoes had the lowest range of variation for the accumulation of sugars,  
200 followed by fresh mid-sized tomatoes. Cherry and cocktail tomatoes represented the highest  
201 levels of variation, not only for sugars but also for glutamic acid. Regarding the rest of acids, mid-  
202 sized fresh tomatoes included higher levels of variation. The general model included a high range  
203 of composition for all the compounds.

204 Regarding the external assay, samples grown in Navarra had similar SSC, but lower sugar and  
205 higher acid contents than those obtained in the sample set 1, corresponding to Extremadura  
206 growing conditions (Table 1).

### 207 ***FT-MIR prediction models***

208 In the three specific sample sets  $R^2$  values for calibration were satisfactory, especially in the case  
209 of sugars and SSC, ranging from 0.82 for fructose in the model with mid-sized tomatoes to 0.98  
210 for SCC in the model of cherry and cocktail tomato (Table 2). The  $R^2$  values for calibration of the  
211 acids tended to be lower, especially in the mid-sized and cherry tomatoes models. For all the  
212 models,  $R^2$  cross-validation values decreased considerably and, at the same time, RMSECV

213 values considerably increased over those obtained for the calibration. Indeed, the cross-  
214 validation seemed to be rather tough, as in the three models the  $R^2$  prediction values were  
215 higher and RMSEP values were similar or lower than RMSECV.

216 In processing tomatoes RMSEP (%mean) values for SSC, citric fructose and glucose were lower  
217 than 10% and slightly higher (<15%) for glutamic and malic acid (Table 2). The error was equally  
218 distributed, and it was similar for samples with low, medium or high contents (fructose  
219 predictions are shown as an example, Fig. 2). Consequently, the relative error tended to  
220 decrease for increasing contents. When RMSEP was contextualized using the maximum value  
221 (%maximum), the percentages of error remained below 9% for all the compounds. In mid-sized  
222 tomatoes, RMSEP (%mean) values were lower than 18% for SSC, glucose, fructose, and citric and  
223 lower than 25% in the rest of compounds. RMSEP (%maximum) values were lower than 10% for  
224 SSC, fructose, glucose and citric and lower than 13% for malic and glutamic acids. In the case of  
225 cherry and cocktail tomatoes RMSEP (%mean) values were lower than 15% for all the  
226 compounds except for glutamic acid (27.5%) and RMSEP (%maximum) values were lower than  
227 10% for all the compounds.

228 Considering the range of variation represented in the samples, for the three models RPD values  
229 were, in general, close to or higher than 2 (Table 2). On the other hand, RER values were higher  
230 than 6 for all the variables in processing and mid-sized-tomatoes, and even higher than 10 for  
231 most variables in cherry and cocktail tomatoes.

232 A general model was obtained with the entire set of samples.  $R^2$  values for calibration ranged  
233 from 0.65 for malic acid to 0.96 for SSC.  $R^2$  and RMSECV values for cross-validation were similar,  
234 to those of the calibration, thus revealing a low impact of cross-validation.  $R^2$  values for the  
235 validation group were almost identical to those of the calibration, except for glutamic acid (0.58  
236 vs. 0.75). RMSEP values even improved RMSEC values (mean decrease of 6%). RMSEP %mean  
237 and %maximum values for main compounds were lower than 15% and 6% respectively. RPD  
238 values were higher than 2.5 for main compounds and RER values were higher than 11 for all the  
239 compounds, except for malic acid (8.5).

240 It seemed that increasing the number of samples did not necessarily lead to a better  
241 performance of the model. In fact, after studying different correlations between model and  
242 prediction parameters and the description of the sample, it resulted that the highest absolute  
243 correlations between relative errors (%mean) and the characteristics of the samples were found  
244 with the coefficient of variation and minimum values of the samples (Supp. Figure 1).  
245 Specifically, the correlation between the coefficients of variation of the sample sets and relative  
246 RMSEC, RMSECV and RMSEP values were 0.75, 0.69 and 0.79 respectively. An inverse  
247 relationship was observed between the minimum value of the samples and RMSEC, RMSECV  
248 and RMSEP values ( $R=-0.46$ ).

249 Sensitivity (true positive rate) and specificity (true negative rate) values were calculated with the  
250 general model to predict the results when a 10 or 20% selection pressure was applied. Sensitivity  
251 for a 10% selection pressure was 100% for SSC, glucose, fructose and citric acid and 75% for  
252 malic and glutamic acid (Table 4). That means that when 10% of samples with the highest  
253 predicted content were selected, all of them had the highest (10%) measured (CZE) values for  
254 the main components. Specificity was 100% for main components and higher than 97% for malic  
255 and glutamic acids, meaning that within the rejected samples only 3% had high malic or glutamic  
256 content. Even for malic and glutamic acid, the mean percentile of CZE measured values of the  
257 selected samples was close to 5%, meaning that false positives also had high contents, close to

258 the limit of selection. When a lower selection pressure was applied (selecting the 20% of  
259 samples), sensitivity decreased, but it was still higher than 78% for main components. Specificity  
260 was higher than 90% in all cases. Mean percentile of CE measured values was below 18% for all  
261 the variables.

## 262 **Prediction of an external assay**

263 The obtained general model was applied to the prediction of a new assay (111 samples) not  
264 included in its development.  $R^2$  values were moderately high in the case of SSC (0.74) and  
265 glucose (0.63) and considerably lower in the rest of the cases (Table 3). Nevertheless, RMSEP  
266 values (%mean) were lower than 15% for SSC and glucose, lower than 20% for fructose and citric  
267 acid and lower than 40% for malic and glutamic acids. RMSEP values (%maximum) were lower  
268 than 21% in all cases. The specific model for sample set 1, representing the same varieties but  
269 grown in different environmental conditions was also used for predictions, but it showed a  
270 worse performance than the general model (Table 3).

271 When a 10% selection pressure was applied to select material using the general model,  
272 sensitivities for main components were higher than 70% and specificities higher than 96% (Table  
273 4). The mean percentile of CZE measured values of the selected samples was lower than 8%.  
274 Thus, the false positives included within the selected samples had high values close to the limit  
275 of selection. With a lower selection pressure (20%) sensitivities decreased, with good values for  
276 SSC and glucose, but lower for fructose (60%) and citric acid (52.4%).

## 277 **DISCUSSION**

278 Tomato is a complex crop, with specific characteristics associated with different varietal types.  
279 For example, it has been proven that the breeding history of processing tomato led to a notable  
280 divergence of the genome from fresh tomato varieties.<sup>27</sup> This divergence is in fact accentuated  
281 in chromosome 5, where several QTL for high soluble solid content can be found. This is in part  
282 due to the strong selection pressure in processing tomato breeding programs for high SSC.  
283 Within fresh tomato, a clear distinction is observed between cherry and other types, such as  
284 round or beef tomatoes, with important differences in the metabolic profile.<sup>28</sup> In this sense,  
285 cherry tomatoes tend to have higher levels of sugars and citric acid and lower levels of malic  
286 acid, as it was observed in the present study.

287 With this level of differentiation, it seems probable that indirect quantification methods for the  
288 accumulation of taste-related compounds should be developed specifically for each varietal  
289 type. In fact, previous works using FT-MIR prediction for sugar and acid content have been either  
290 focussed on processing tomato,<sup>19,20</sup> or using a limited number of varieties. That is the case of  
291 Beullens et al.<sup>29</sup> that applied their model to four varieties or Vermeir et al.<sup>25</sup> with sets of six and  
292 three varieties. Nevertheless, some attempts have been made to include a higher level of  
293 diversity. It is the case of the study performed by Scibisz et al.<sup>18</sup>, that included 39 commercial  
294 and traditional varieties of, i.a., cherry, cocktail, Marmande, and processing tomato.

295 Consequently, the results obtained in these studies are not easily compared. Usually, cross-  
296 validation performance is provided to assess the robustness of the models. In the present study  
297 cross-validation using Venetian blinds proved to be rather tough, as the values for  $R^2$  and RMSE  
298 in the prediction improved those of the cross-validation. Nevertheless, the performance of the  
299 model predicting values for samples not included in the calibration gives a better idea of its  
300 robustness, thus comparisons should be made with prediction performance results. But even  
301 then, not always mean values and coefficients of variation are provided to contextualize RMSEP



302 values, and other values such as the dimensionless RPD and RER values are not usually provided.  
303 RPD can be interpreted as the ratio of natural variation in the samples to the size of likely  
304 prediction errors and RER is similar in spirit, but in this case, it uses the range of variation of the  
305 samples instead of the standard deviation.<sup>24</sup>

306 Despite these limitations, some comparisons can be made. In the case of processing tomato,  
307 Ayvaz et al.<sup>20</sup> obtained RMSEP values for glucose of 1.4 g L<sup>-1</sup> and fructose 1.46 g L<sup>-1</sup> and SSC of  
308 0.12°Brix, values that contextualized with the mean result in 16%, 17% and 2% for ranges 1.1-  
309 20.2 g L<sup>-1</sup>, 0.1-20.4 g L<sup>-1</sup> and 3.8-7.2°Brix respectively. Also with processing tomato, Wilkerson et  
310 al.<sup>19</sup> obtained values of RMSEP with a triple-bounce ATR of 1.47 g L<sup>-1</sup>, 1.23 g L<sup>-1</sup> and 0.23 g L<sup>-1</sup> for  
311 glucose, fructose and SSC. In this case, the means were not provided, but samples used for  
312 prediction ranged 10.0-21.4 g L<sup>-1</sup> for glucose, 11-20.6 g L<sup>-1</sup> for fructose and 4.2-6.7 °Brix for SSC.  
313 These values are similar to those obtained in the present work in the model restricted to  
314 processing tomato (RMSEP values of 1.2 g kg<sup>-1</sup>, 1.4 g kg<sup>-1</sup> and 0.2 °Brix with ranges of 10.01-19.69  
315 g kg<sup>-1</sup>, 9.03-20.04 g kg<sup>-1</sup> and 3.5-5.8 °Brix respectively).

316 Regarding studies with a wider diversity, Scibisz et al.<sup>18</sup> obtained RMSEP values (contextualized  
317 with the mean) for glucose, fructose and SSC of 5.1%, 6.8%, and 2.9% respectively with  
318 coefficients of variation of prediction samples of 25%, 20% and 20% respectively. While in our  
319 case with the general model RMSEP values were 12.1%, 14.3% and 4.9% with much higher  
320 variability in the prediction samples (coefficients of variation of 43.2%, 39.9%, and 26.7%  
321 respectively).

322 Few studies are available where the authors have applied a PLS regression model to samples  
323 obtained in an external assay. Among them, Scibisz et al.<sup>18</sup> obtained models for individual  
324 compounds, but the external assay was applied only to predict SSC, total acidity and dry matter.  
325 In that case, they obtained an RMSEP (%mean) value for SSC of 4.5%. In the present work, RMSEP  
326 was higher 8.4%, but it should be considered that Scibisz et al.<sup>18</sup> included two varieties and the  
327 present work included 8 varieties grown with different water and fertilization regimes.<sup>21,22</sup> It  
328 could be argued that the external sample set included varieties already considered in the model  
329 and that the samples were obtained in the same years (those studies were performed with  
330 samples obtained during two years). It should be considered though, that the general model  
331 included not only these eight varieties, but much more. Additionally, previous studies with these  
332 materials proved that the site of cultivation had a similar effect or even higher than the year  
333 effect on the metabolic profile of tomato.<sup>21,22</sup> Even considering the limitation of the external  
334 sample set, it seems clear that the potential of using FT-MIR models with external samples would  
335 be a real.

336 The models developed in this work have proved to be robust and reliable for main components  
337 (SSC, glucose, fructose and citric acid), with worse models obtained for malic and glutamic acid.  
338 Vermeir et al.<sup>25</sup> also found lower predictivity levels for malic and glutamic acids compared to  
339 citric acid. In that case, they justified that the range of variation for these compounds and their  
340 concentration was small. Scibisz et al.<sup>18</sup> also found that the model for malic acid was weak, and  
341 also stated that it was probably due to the low concentration found in tomato, in agreement  
342 with the results of Rudnitskaya et al.<sup>30</sup> in apple, pointing out that FT-MIR models were not  
343 suitable for minor components. On the other hand, Wilkerson et al.<sup>19</sup> did not include malic acid  
344 in their models, but the characteristics of the glutamic model were similar to those of the citric  
345 acid. On the other hand, most published literature shows that models for SSC are much better  
346 than those for individual components<sup>18-20</sup>, as in the present work. One plausible explanation is  
347 that SSC involves the signal of individual components and it is obviously higher. Nevertheless,

348 SSC indirect prediction has the only advantage of providing an estimate to contextualize new  
349 results with previous experiences were no data for individual compounds is available.

350 Apart from different comparisons, it seems clear that the obtained general model will be highly  
351 valuable for screening purposes. RPD values higher than 2 and RER values higher than 10,  
352 confirm this point, as these values are considered a threshold to define useful models.<sup>24-26</sup> In  
353 addition, RMSEP (%maximum) values for SSC, fructose, glucose and citric acid were lower than  
354 6%. This evidences the goodness of the indirect FT-MIR determinations in the selection of  
355 materials with outstanding contents of taste-related compounds.

356 It also seems clear that a general model is most robust than specific models for each varietal  
357 type. In our case, this seems to be due to the higher variability included rather than to the use  
358 of a higher number of samples, as the performance of the models represented in the RMSE  
359 (%mean) values, had a high positive correlation with the coefficient of variation of the sample  
360 sets used during the calibration. It also seems that low minimum values of the variables in the  
361 samples would interfere with the calculation of good models, as moderate negative correlations  
362 were obtained with RMSE (%mean) values. Thus, answering the initial question proposed in this  
363 work, a higher amplitude of samples and environments really increases the efficiency of FT-MIR  
364 models in the prediction of the concentration of taste-related compounds in tomato.

365 The feasibility of the use of FT-MIR models in screening activities during breeding programs is  
366 clear. The benefits of the use of indirect quantification methods are evident when it is  
367 considered that in screening programs for a high content of specific sugars and acids, a notable  
368 contribution of genotype x environment interaction has been described, forcing to develop  
369 multi-environmental trials.<sup>31</sup> Even more, the considerable amount of variation present within  
370 accessions, especially in wild germplasm, leads to the need to analyze individually a high number  
371 of plants per accession in order to identify sources of variation. And again, a high number of  
372 plants are to be analyzed in segregant populations during the introgression of the trait.

373 The sensitivity of the PLS regression model when applying a 10% selection pressure is 100%.  
374 That means that in order to select the best materials in screening programs, time-consuming  
375 and expensive HPLC or CZE determinations can be replaced by rapid and cheap ATR FT-MIR  
376 analysis. This is of utmost utility for the development of breeding programs, which still rely on  
377 the calculation of gross measures such as SSC due to the cost of implementing selections based  
378 on individual compounds.

379 In indirect spectroscopic analysis, it is usually necessary to perform a calibration for each assay.  
380 But in the present work, it has been proven that even applying the general model to an external  
381 assay not included in the calibration, as it was possible to obtain sensitivities higher than 70%  
382 with 10% selection pressure for main taste-related compounds. Additionally, those false  
383 positives corresponded to samples with high contents close to the limit of selection, as the mean  
384 percentiles of selected samples were close to the ideal 5%. Thus, even if these samples were  
385 selected, their effect on the development of a breeding program would be minimum. Even for  
386 quantification purposes, RMSEP values contextualized with maximum values were lower than  
387 13% for main compounds, which is a notable performance, considering that it has been obtained  
388 with samples not included in calibration models.

389

390 **CONCLUSIONS**

391 Our results support the applicability of ATR FT-MIR as a tool to select samples with high contents  
392 of SSC, glucose, fructose and citric acid. For that purpose, it would be advisable to use pools of  
393 samples with the highest variability to develop robust general models, rather than specific  
394 models. Preliminary results using an external sample set with materials grown in a different  
395 environment suggest that even though a calibration for each assay is advisable, it would not be  
396 necessary for screening programs if a later more precise determination is planned. Even if no  
397 further analysis were performed, most selected samples would be close to the selection  
398 objective.

399

#### 400 **ACKNOWLEDGMENTS**

401 This research was performed despite the lack of direct public funding for its development and  
402 thanks to the enthusiasm of the authors. The authors thank Dr. Lahoz and Dr. Campillo for  
403 providing samples of processing tomato. G. Ibañez thanks Universitat Jaume I for funding his  
404 pre-doctoral grant (PREDOC/2015/45).

405

#### 406 **REFERENCES**

407

- 408 1. Bruhn CM, Feldman N, Garlitz C, Harwood J, Ivans E, Marshall M, et al. Consumer  
409 perceptions of quality: apricots, cantaloupes, peaches, pears, strawberries,. *J Food Qual*  
410 **14**(3):187–195 (1991).<https://doi.org/10.1111/j.1745-4557.1991.tb00060.x>.
- 411 2. Cebolla-Cornejo J, Soler S, Nuez F. Genetic erosion of traditional varieties of vegetable  
412 crops in Europe: tomato cultivation in Valencia (Spain) as a case Study. *Int J Plant Prod*  
413 **1**(2):113–128 (2007).<https://doi.org/10.22069/IJPP.2012.531>.
- 414 3. Davies JN, Hobson GE. The constituents of tomato fruit — the influence of environment,  
415 nutrition, and genotype. *C R C Crit Rev Food Sci Nutr* **15**(3):205–280  
416 (1981).<https://doi.org/10.1080/10408398109527317>.
- 417 4. Díaz de León-Sánchez F, Pelayo-Zaldívar C, Rivera-Cabrera F, Ponce-Valadez M,  
418 Ávila-Alejandre X, Fernández FJ, et al. Effect of refrigerated storage on aroma and  
419 alcohol dehydrogenase activity in tomato fruit. *Postharvest Biol Technol* **54**(2):93–100  
420 (2009).<https://doi.org/10.1016/j.postharvbio.2009.07.003>.
- 421 5. Cebolla-Cornejo J, Roselló S, Valcárcel M, Serrano E, Beltrán J, Nuez F. Evaluation of  
422 Genotype and Environment Effects on Taste and Aroma Flavor Components of Spanish  
423 Fresh Tomato Varieties. *J Agric Food Chem* **59**(6):2440–2450  
424 (2011).<https://doi.org/10.1021/jf1045427>.
- 425 6. Tieman D, Zhu G, Resende MFR, Lin T, Nguyen C, Bies D, et al. A chemical genetic  
426 roadmap to improved tomato flavor. *Science (80- )* **355**:391–394 (2017).
- 427 7. Powell ALT, Nguyen CV, Hill T, Cheng KL, Figueroa-Balderas R, Aktas H, et al. Uniform  
428 ripening Encodes a Golden. *Science (80- )* **336**(29):1711–1715  
429 (2012).<https://doi.org/10.1126/science.1222218>.
- 430 8. Osorio S, Alba R, Damasceno CM, Lopez-Casado G, Lohse M, Zanon MI, et al. Systems  
431 Biology of Tomato Fruit Development: Combined Transcript, Protein, and Metabolite  
432 Analysis of Tomato Transcription Factor (nor, rin) and Ethylene Receptor (Nr) Mutants  
433 Reveals Novel Regulatory Interactions. *Plant Physiol* **157**(1):405–425  
434 (2011).<https://doi.org/10.1104/pp.111.175463>.
- 435 9. Qin G, Wang Y, Cao B, Wang W, Tian S. Unraveling the regulatory network of the  
436 MADS box transcription factor RIN in fruit ripening. *Plant J* **70**(2):243–255

- 437 (2012).<https://doi.org/10.1111/j.1365-313X.2011.04861.x>.
- 438 10. Baldwin EA, Scott JW, Einstein M, Malundo TMM, Carr BT, Shewfelt RL, et al.  
439 Relationship between sensory and instrumental analysis for tomato flavor. *J Am Soc*  
440 *Hortic Sci* **123**(5):906–915 (1998).
- 441 11. Causse M, Friguet C, Coiret C, L  picier M, Navez B, Lee M, et al. Consumer  
442 Preferences for Fresh Tomato at the European Scale: A Common Segmentation on  
443 Taste and Firmness. *J Food Sci* **75**(9):531–541 (2010).<https://doi.org/10.1111/j.1750-3841.2010.01841.x>.
- 445 12. Agius C, von Tucher S, Poppenberger B, Rozhon W. Quantification of sugars and  
446 organic acids in tomato fruits. *MethodsX* **5**(May):537–550  
447 (2018).<https://doi.org/10.1016/j.mex.2018.05.014>.
- 448 13. Cebolla-Cornejo J, Valc  rcel M, Herrero-Mart  nez JM, Rosell   S, Nuez F. High efficiency  
449 joint CZE determination of sugars and acids in vegetables and fruits. *Electrophoresis*  
450 **33**(15):2416–2423 (2012).<https://doi.org/10.1002/elps.201100640>.
- 451 14. Bureau S, Cozzolino D, Clark CJ. Contributions of Fourier-transform mid infrared (FT-  
452 MIR) spectroscopy to the study of fruit and vegetables: A review. *Postharvest Biol*  
453 *Technol* **148**:1–14 (2019).<https://doi.org/10.1016/j.postharvbio.2018.10.003>.
- 454 15. Nicolai BM, Beullens K, Bobelyn E, Peirs A, Saeys W, Theron KI, et al. Nondestructive  
455 measurement of fruit and vegetable quality by means of NIR spectroscopy: A review.  
456 *Postharvest Biol Technol* **46**(2):99–118  
457 (2007).<https://doi.org/10.1016/j.postharvbio.2007.06.024>.
- 458 16. Garc  a-Mart  nez, S., G  lvez-Sola, L. N., Alonso, A., Agull  , E., Rubio, F., Ruiz, J. J., &  
459 Moral, R. Quality assessment of tomato landraces and virus-resistant breeding lines:  
460 quick estimation by near infrared reflectance spectroscopy. *Journal of the Science of*  
461 *Food and Agriculture*, **92**(6), 1178-1185 (2012). <https://doi.org/10.1002/jsfa.4661>.
- 462 17. de Oliveira GA, de Castilhos F, Renard CMGC, Bureau S. Comparison of NIR and MIR  
463 spectroscopic methods for determination of individual sugars, organic acids and  
464 carotenoids in passion fruit. *Food Res Int* **60**:154–162  
465 (2014).<https://doi.org/10.1016/j.foodres.2013.10.051>.
- 466 18. Scibisz I, Reich M, Bureau S, Gouble B, Causse M, Bertrand D, et al. Mid-infrared  
467 spectroscopy as a tool for rapid determination of internal quality parameters in tomato.  
468 *Food Chem* **125**(4):1390–1397 (2011).<https://doi.org/10.1016/j.foodchem.2010.10.012>.
- 469 19. Wilkerson ED, Anthon GE, Barrett DM, Sayajon GFG, Santos AM, Rodriguez-Saona LE.  
470 Rapid assessment of quality parameters in processing tomatoes using hand-held and  
471 benchtop infrared spectrometers and multivariate analysis. *J Agric Food Chem*  
472 **61**(9):2088–2095 (2013).<https://doi.org/10.1021/jf304968f>.
- 473 20. Ayvaz H, Sierra-Cadavid A, Aykas DP, Mulqueeney B, Sullivan S, Rodriguez-Saona LE.  
474 Monitoring multicomponent quality traits in tomato juice using portable mid-infrared (MIR)  
475 spectroscopy and multivariate analysis. *Food Control* **66**:79–86  
476 (2016).<https://doi.org/http://dx.doi.org/10.1016/j.foodcont.2016.01.031>.
- 477 21. Lahoz I, Leiva-Brondo M, Mart   R, Macua JI, Campillo C, Rosell   S, et al. Influence of  
478 high lycopene varieties and organic farming on the production and quality of processing  
479 tomato. *Sci Hortic (Amsterdam)* **204**:128–137  
480 (2016).<https://doi.org/10.1016/j.scienta.2016.03.042>.
- 481 22. Mart   R, Valc  rcel M, Leiva-Brondo M, Lahoz I, Campillo C, Rosell   S, et al. Influence of  
482 controlled deficit irrigation on tomato functional value. *Food Chem* **252**:250–257  
483 (2018).<https://doi.org/10.1016/j.foodchem.2018.01.098>.
- 484 23. Naes T, Isaksson T, Fearn T, Davies T. A User-Friendly Guide to Multivariate Calibration  
485 and Classification. NIR Publications; 2002.
- 486 24. Fearn T. Assessing calibrations: SEP, RPD, RER and R2. *NIR news* **13**(6):12–14

- 487 (2002).<https://doi.org/10.1255/nirn.689>.
- 488 25. Vermeir S, Beullens K, Mészáros P, Polshin E, Nicolai BM, Lammertyn J. Sequential  
489 injection ATR-FTIR spectroscopy for taste analysis in tomato. *Sensors Actuators, B*  
490 *Chem* **137**(2):715–721 (2009).<https://doi.org/10.1016/j.snb.2009.01.056>.
- 491 26. Williams P, Norris K. Near-infrared technology in the agricultural and food industries.  
492 American Association of Cereal Chemists, Inc.; 1987.
- 493 27. Lin T, Zhu G, Zhang J, Xu X, Yu Q, Zheng Z, et al. Genomic analyses provide insights  
494 into the history of tomato breeding. *Nat Genet* **46**(11):1220–1226  
495 (2014).<https://doi.org/10.1038/ng.3117>.
- 496 28. Ursem R, Tikunov Y, Bovy A, Van Berloo R, Van Eeuwijk F. A correlation network  
497 approach to metabolic data analysis for tomato fruits. *Euphytica* **161**(1–2):181–193  
498 (2008).<https://doi.org/10.1007/s10681-008-9672-y>.
- 499 29. Beullens K, Kirsanov D, Irudayaraj J, Rudnitskaya A, Legin A, Nicolai BM, et al. The  
500 electronic tongue and ATR-FTIR for rapid detection of sugars and acids in tomatoes.  
501 *Sensors Actuators, B Chem* **116**(1–2):107–115  
502 (2006).<https://doi.org/10.1016/j.snb.2005.11.084>.
- 503 30. Rudnitskaya A, Kirsanov D, Legin A, Beullens K, Lammertyn J, Nicolai BM, et al.  
504 Analysis of apples varieties – comparison of electronic tongue with different analytical  
505 techniques. *Sensors Actuators B Chem* **116**(1–2):23–28  
506 (2006).<https://doi.org/10.1016/j.snb.2005.11.069>.
- 507 31. Galiana-Balaguer L, Ibáñez G, Cebolla-Cornejo J, Roselló S. Evaluation of germplasm in  
508 *Solanum* section *Lycopersicon* for tomato taste improvement. *Turkish J Agric For*  
509 **42**(5):309–321 (2018).<https://doi.org/10.3906/tar-1712-61>.

510 **Table 1.** Statistical parameters of the samples regarding the accumulation of soluble solids content (SSC), sugars and acids contents using refractometry and  
 511 capillary zone electrophoresis. For each sample set the characteristics of the samples used for calibration ( $n_c$ ) and cross-validation ( $n_{cv}$ ), and those used for  
 512 prediction are indicated. An external assay was included only to predict values using the general model.

	Parameters	Calibration group samples					Validation group samples				
		Mean	Standard deviation	Coefficient of variation (%)	Minimum	Maximum	Mean	Standard deviation	Coefficient of variation (%)	Minimum	Maximum
Sample set 1 Processing tomato Extremadura ( $n_c = 81$ ; $n_v = 27$ )	SSC <sup>1</sup>	4.55	0.42	9.2	3.80	5.55	4.54	0.60	13.1	3.50	5.80
	Glucose <sup>2</sup>	14.16	2.08	14.7	9.84	19.48	14.57	2.50	17.2	10.01	19.69
	Fructose <sup>2</sup>	14.83	2.46	16.6	9.72	19.94	15.14	2.97	19.6	9.03	20.04
	Citric <sup>2</sup>	3.53	0.57	16.0	2.25	5.34	3.62	0.60	16.7	2.53	5.22
	Malic <sup>2</sup>	1.23	0.30	24.1	0.79	1.83	1.25	0.30	23.7	0.87	2.00
	Glutamic <sup>2</sup>	0.96	0.43	44.7	0.36	2.35	1.03	0.41	39.7	0.42	1.96
Sample set 2 Fresh mid-sized tomato ( $n_c = 80$ ; $n_v = 27$ )	SSC <sup>1</sup>	4.27	0.55	13.0	2.95	5.40	4.25	0.52	12.3	3.50	5.65
	Glucose <sup>2</sup>	12.10	3.40	28.0	6.73	22.49	11.50	3.31	28.8	6.98	22.72
	Fructose <sup>2</sup>	13.77	3.67	26.7	8.25	25.63	12.99	3.47	26.7	8.15	24.18
	Citric <sup>2</sup>	5.84	1.93	33.0	3.06	14.03	5.43	1.74	32.1	2.70	9.75
	Malic <sup>2</sup>	1.76	0.63	35.6	0.56	4.02	1.73	0.69	39.7	0.56	3.75
	Glutamic <sup>2</sup>	1.70	0.76	44.8	0.71	4.25	1.68	0.68	40.3	0.70	3.25
Sample set 3 Cherry&cocktail tomato ( $n_c = 86$ ; $n_v = 29$ )	SSC <sup>1</sup>	5.77	1.44	24.9	3.50	11.35	6.00	1.84	30.7	4.00	12.40
	Glucose <sup>2</sup>	18.14	7.73	42.6	6.73	46.09	19.33	9.13	47.2	10.03	50.42
	Fructose <sup>2</sup>	20.81	7.77	37.3	8.44	51.99	21.57	9.15	42.4	12.41	53.43
	Citric <sup>2</sup>	8.42	1.65	19.6	4.67	12.38	8.36	1.60	19.1	5.20	11.66
	Malic <sup>2</sup>	1.47	0.47	31.6	0.79	3.44	1.50	0.49	32.9	0.96	3.01
	Glutamic <sup>2</sup>	1.60	1.18	73.4	0.32	6.77	1.71	1.47	86.1	0.28	6.43
General model ( $n_c = 245$ ; $n_v = 87$ )	SSC <sup>1</sup>	4.91	1.18	23.9	2.95	11.35	4.90	1.31	26.7	3.50	12.40
	Glucose <sup>2</sup>	14.97	5.70	38.1	6.73	46.09	15.17	6.55	43.2	6.73	43.18
	Fructose <sup>2</sup>	16.66	6.17	37.0	8.15	51.99	16.71	6.67	39.9	8.25	53.43
	Citric <sup>2</sup>	5.90	2.43	41.2	2.36	12.38	6.13	2.58	42.1	2.25	14.03
	Malic <sup>2</sup>	1.51	0.54	35.9	0.56	4.02	1.45	0.47	32.4	0.56	3.01
	Glutamic <sup>2</sup>	1.47	0.93	63.1	0.28	6.77	1.43	0.92	64.1	0.32	5.12
External assay Processing tomato Navarra ( $n_v = 111$ )	SSC <sup>1</sup>						4.66	0.58	12.5	3.45	6.10
	Glucose <sup>2</sup>						12.38	2.82	22.8	6.10	20.87
	Fructose <sup>2</sup>						13.13	3.14	24.6	5.84	22.42
	Citric <sup>2</sup>						4.21	0.85	20.2	2.13	7.06
	Malic <sup>2</sup>						0.84	0.24	28.1	0.32	1.49
	Glutamic <sup>2</sup>						1.70	0.42	24.6	0.81	2.82

513 <sup>1</sup>°Brix; <sup>2</sup> g kg<sup>-1</sup>

514 **Table 2.** Performance of the partial least squares (PLS) regression models predicting taste-related compounds content from ATR FT-MIR spectra. SSC: soluble  
 515 solids content; R<sup>2</sup> correlation coefficient; RMSE: root mean squared error; C: calibration; CV: cross-validation; P: prediction; RPD: residual prediction deviation;  
 516 RER: range error ratio.

Model	Parameters	R <sup>2</sup> <sub>C</sub>	RMSEC	R <sup>2</sup> <sub>CV</sub>	RMSECV	R <sup>2</sup> <sub>P</sub>	RMSEP	RMSEP (%Maximum)	RMSEP (%Mean)	RPD	RER
Sample set 1 Processing tomato Extremadura (n <sub>c</sub> <sup>3</sup> =81; n <sub>v</sub> <sup>4</sup> = 27)	SSC <sup>1</sup>	0.85	0.2	0.50	0.3	0.81	0.2	4.2	5.3	2.5	9.5
	Glucose <sup>2</sup>	0.88	0.7	0.63	1.3	0.79	1.2	5.9	8.0	2.2	8.3
	Fructose <sup>2</sup>	0.88	0.9	0.45	1.9	0.77	1.4	6.8	9.0	2.2	8.1
	Citric <sup>2</sup>	0.90	0.2	0.51	0.4	0.68	0.3	5.6	8.0	2.1	9.3
	Malic <sup>2</sup>	0.82	0.1	0.52	0.2	0.58	0.2	8.5	13.6	1.8	6.7
	Glutamic <sup>2</sup>	0.94	0.1	0.70	0.2	0.83	0.1	7.1	13.6	2.9	11.0
Sample set 2 Fresh mid-sized tomato (n <sub>c</sub> =80; n <sub>v</sub> = 27)	SSC <sup>1</sup>	0.93	0.1	0.78	0.3	0.91	0.2	3.3	4.4	2.8	11.6
	Glucose <sup>2</sup>	0.83	1.4	0.39	2.7	0.57	2.0	8.8	17.4	1.7	7.9
	Fructose <sup>2</sup>	0.82	1.5	0.30	3.1	0.70	1.6	6.7	12.4	2.2	10.0
	Citric <sup>2</sup>	0.58	1.1	0.39	1.3	0.75	0.9	9.6	17.1	1.9	7.6
	Malic <sup>2</sup>	0.63	0.4	0.09	0.6	0.49	0.4	10.7	23.1	1.7	8.0
	Glutamic <sup>2</sup>	0.54	0.5	0.38	0.6	0.65	0.4	12.6	24.4	1.7	6.2
Sample set 3 Cherry & cocktail tomato (n <sub>c</sub> =86; n <sub>v</sub> = 29)	SSC <sup>1</sup>	0.98	0.2	0.97	0.3	0.98	0.3	2.6	5.3	5.8	26.5
	Glucose <sup>2</sup>	0.99	0.9	0.94	1.9	0.99	1.1	2.1	5.4	8.7	38.5
	Fructose <sup>2</sup>	0.95	1.8	0.92	2.2	0.93	2.4	4.5	11.2	3.8	17.0
	Citric <sup>2</sup>	0.76	0.8	0.60	1.0	0.57	1.1	9.0	12.6	1.5	6.1
	Malic <sup>2</sup>	0.92	0.1	0.62	0.3	0.90	0.2	5.3	10.7	3.1	12.8
	Glutamic <sup>2</sup>	0.88	0.4	0.82	0.5	0.89	0.5	7.3	27.5	3.1	13.1
General model (n <sub>c</sub> =245; n <sub>v</sub> = 87)	SSC <sup>1</sup>	0.96	0.2	0.94	0.3	0.95	0.2	1.9	4.9	5.5	37.1
	Glucose <sup>2</sup>	0.92	1.6	0.89	1.8	0.89	1.8	4.3	12.1	3.6	19.8
	Fructose <sup>2</sup>	0.89	2.1	0.84	2.5	0.87	2.4	4.5	14.3	2.8	18.9
	Citric <sup>2</sup>	0.95	0.6	0.87	0.9	0.90	0.8	5.6	12.7	3.3	15.1
	Malic <sup>2</sup>	0.65	0.3	0.54	0.3	0.54	0.3	9.6	20.0	1.6	8.5
	Glutamic <sup>2</sup>	0.85	0.4	0.75	0.5	0.58	0.4	8.0	28.7	2.2	11.7

517 <sup>1</sup>RMSE values expressed as °Brix; <sup>2</sup>RMSE values expressed as g kg<sup>-1</sup>; <sup>3</sup>number of samples in calibration set; <sup>4</sup>number of samples in validation set

518

519

**Table 3.** Performance of the predictions made for the external assay (processing tomato grown in Navarra; 111 samples) obtained with the general model and the specific model of sample set 1 (processing tomato grown in Extremadura).  $R^2_p$  correlation coefficient of the prediction; RMSEP: root mean squared error of the prediction ( $^{\circ}$ Brix for SSC and  $\text{g kg}^{-1}$  for individual compounds).

Model	Parameter	$R^2_p$	RMSEP	RMSEP (%Mean)	RMSEP (%Maximum)
General model predicting external assay (processing tomato Navarra)	SSC	0.74	0.39	8.4	6.4
	Glucose	0.63	1.74	14.1	8.3
	Fructose	0.34	2.51	19.1	11.2
	Citric	0.25	0.86	20.4	12.2
	Malic	0.18	0.31	36.9	20.8
	Glutamic	0.19	0.49	28.6	17.2
Sample set 1 (processing tomato Extremadura) model predicting external assay (processing tomato Navarra)	SSC	0.56	0.68	14.6	11.12
	Glucose	0.49	2.40	19.4	11.5
	Fructose	0.17	3.66	27.9	16.3
	Citric	0.05	1.07	25.4	15.2
	Malic	0.21	0.24	28.8	16.23
	Glutamic	0.11	0.45	26.2	15.8



**Table 4.** Values of sensitivity and specificity obtained using the general PLS model in plant selection for high content of taste-related constituents when applying a selection pressure of 10% or 20%. Samples from the general model validation set group and the external assay (processing tomato grown in Navarra) were evaluated. Mean percentiles of selected plants are also provided.

General model	10% Selection pressure			20% Selection pressure		
	Sensitivity (%)	Specificity (%)	Mean percentile (%)	Sensitivity (%)	Specificity (%)	Mean percentile (%)
SSC	100.0	100.0	5.8	86.7	96.8	11.3
Glucose	100.0	100.0	5.6	87.5	96.9	11.7
Fructose	100.0	100.0	5.4	82.4	92.4	16.0
Citric	100.0	100.0	5.8	78.6	94.6	12.3
Malic	75.0	97.1	5.3	62.5	90.3	13.3
Glutamic	75.0	97.2	6.9	68.8	92.2	17.7

External assay	10% Selection pressure			20% Selection pressure		
	Sensitivity (%)	Specificity (%)	Mean percentile (%)	Sensitivity (%)	Specificity (%)	Mean percentile (%)
SSC	88.9	98.8	5.4	73.7	93.4	13.2
Glucose	90.0	98.9	5.8	70.0	92.5	14.0
Fructose	70.0	96.7	7.0	60.0	90.1	21.9
Citric	70.0	96.8	8.0	52.4	88.0	22.0
Malic	40.0	93.6	8.8	47.6	86.6	19.3
Glutamic	60.0	95.7	10.0	42.9	85.4	37.1

## FIGURE CAPTIONS

**Figure 1.** FT-MIR spectra of tomato samples in the 4000-650  $\text{cm}^{-1}$  region (a) and in the 1500-900  $\text{cm}^{-1}$  region (b) used for chemometric analysis.

**Figure 2.** Predicted (ATR FT-MIR) vs. measured (CZE) glucose contents using the partial least squares (PLS) regression model for sample set 1 (processing tomato in Extremadura). Grey dots: samples used to calibrate the PLS model; red diamonds: samples used to predict using the PLS model.