# RSC Advances

ROYAL SOCIETY OF CHEMISTRY

www.rsc.org/advances

**RSC Advances**

# Graphical abstract

# Discrimination of Moldy Wheat Using Terahertz Imaging Combined with Multivariate Classification

Yuying Jiang,[a,c] Hongyi Ge, [b*] Feiyu Lian,[b] Yuan Zhang,[b] and Shanhong Xia[a]

[a]*State Key Laboratory of Transducer Technology, Institute of Electronics, Chinese Academy of Sciences, Beijing 100080, China*

[b]*Key Laboratory of Grain Information Processing& Control, Ministry of Education, Henan University of Technology, Zhengzhou 450001, China*

[c]*University of Chinese Academy of Sciences, Beijing 100080, China*

*\* Email：gehongyi@haut.edu.cn, phone: 0086-371-67756610*

Terahertz (THz) imaging was employed to develop a novel method for discriminating wheat of varying states of moldiness. Spectral data, in the range of 0.2–1.6 THz, were extracted from regions of interest (ROIs) in the THz images. Principal component analysis (PCA) was used to evaluate the spectral data and determine the cluster trend. Six optimal frequencies were selected by implementing PCA directly for each image's ROI. Classification models for moldy wheat identification were established using the support vector machine (SVM) method, a partial least-squares regression analysis, and the back propagation neural network method. The models developed from these methods were based on the full and optimal frequencies, using the top three principal components as input variables. The PCA-SVM method achieved a prediction accuracy of over 95%, and was implemented at every pixel in the images to visually demonstrate the moldy wheat classification method. Our results indicate that THz imaging combined with chemometric algorithms is efficient and practical for the discrimination of moldy wheat.

26 **Keywords: discriminate analysis, wheat grain, moldy, terahertz imaging, spectral analysis,**

27 **multivariate classification**

28

## 1. Introduction

Wheat is a primary food crop worldwide, and contains high amounts of carbohydrates, proteins, fat, and vitamins (Oladunmoye, Akinoso, & Olapade, 2010). Mildew such as Aflatoxinand and Aspergillusniger are prevalent throughout all stages of wheat growth and production. When improperly stored and processed, these mildews pose a potential threat to humans and fowls (Neethirajan, Karunakaran, Jayas, & White, 2007). Recently, food quality and safety assessment have increased within the food industry. Conventional moldy grain detection methods, such as naked-eye observations, microscope inspection, liquid chromatography, and enzyme-linked immunosorbent assays, are time-consuming and labor-intensive (Turner, Subrahmanyam, & Piletsky, 2009).

To satisfy the demand for high-quality consumer products, extensive studies into grain quality via nondestructive rapid evaluations have been performed. Wang et al. (Wang, Liu, Yu, Wu & He, 2011) presented a new approach for non-invasive classification of raisins by using computer vision techniques. Eifler et al. (Eifler, Martinelli, Santonico, Capuano, Schild, & Di Natale, 2011) used an electronic nose to differentiate between infected and non-infected wheat grains. Arngren et al. (Arngren, Hansen, Eriksen, Larsen, &Larsen, 2011) used near-infrared hyperspectral imaging combined with nonlinear neural networks to identify early-stage pregermination in barley grains. ElMasry et al. (ElMasry, Wang, ElSayed, Ngadi, 2007) proposed a novel tool for nondestructive determination of moisture content, total soluble solids, and acidity in strawberry using NIR spectroscopy. However, these measurement techniques do not probe the far-infrared spectral region, which contains a wealth of physical

50    and chemical information.

51    Terahertz (THz) radiation (with frequencies from 0.3 to 10 THz and wavelengths from 3.3

52    to 333 cm$^{-1}$ ) occupies the region between the microwave and infrared bands; it can be used for

53    non-destructive and non-invasive analyses, and possesses attractive features such as extremely

54    low-energy levels, broad spectral bandwidth, transparency, and good penetration through

55    various materials (Ferguson, & Zhang, 2002; Tonouchi, 2007). THz spectroscopy and imaging

56    are rapidly becoming novel techniques in the field of optics research. The new techniques are

57    widely used as solutions in art conservation (Fukunaga, & Hosako, 2010), security problems

58    (Melinger, Laman, & Grischkowsky, 2008), biomedical applications (Oh. et. al., 2014; Siegel,

59    2004), agricultural quality control (Gowen, O'Sullivan, & O'Donnell, 2012; Ge, Jiang, Xu, Lian,

60    Zhang, & Xia, 2014), and other fields (Guillet. et. al., 2014). THz imaging is performed both by

61    the transmission and reflection of THz waves. In reflectance imaging, THz waves reflect not

62    only from the surface of samples, but also from interfaces present in the samples within the

63    penetration depth of the radiation (Safrai, Ben Ishai, Polsman, Einav, & Feldman, 2014). Thus,

64    both surface and depth information can be obtained from the timing and amplitude of the

65    reflected waves. Time- and frequency-domain structural images can be acquired from detected

66    THz waves associated with various parameters at each pixel in the measured sample area (Reid,

67    Pickwell-MacPherson, Laufer, Gibson, Hebden, & Wallace, 2010). Owing to the absorption,

68    reflection, scattering, and phase-shifting of the imaged material, measured parameters can

69    change due to differing wave delay and attenuation.

70    The aim of this study was to evaluate the validity and feasibility of identifying different

71  moldy states of wheat using THz imaging and multivariate data analysis methods. THz spectra

72  of wheat grains with different moldy statuses were extracted in the range of 0.2–1.6 THz from

73  regions of interest (ROIs) in each THz image. Principal component analysis (PCA) was used to

74  explore features of the spectral data and select the optimal frequencies. Support vector machine

75  (SVM), partial least-squares regression (PLSR), and back propagation neural network (BPNN)

76  models were established based on the full frequencies and optimal frequencies for

77  discriminating between the four stages of moldy wheat. Finally, THz images of wheat with

78  different moldy states were investigated using the optimal classification method (i.e.,

79  PCA-SVM).

80  **2. Materials and methods**

81  *2.1 Experimental setup*

82  A standard THz-TDS laboratory setup, using reflection geometry as developed by Zomega

83  Terahertz Corporation in USA, was used in our experiment. A schematic of the THz-TDS

84  reflection imaging system is shown in Fig. 1. The THz imaging system employed an externally

85  pulsed femtosecond laser Ti-sapphire with a pulse width, central wavelength, and repetition

86  frequency of 100 fs, 800 nm, and 80 MHz, respectively. The beam produced by the laser was

87  split into a pump and a probe using a polarizing beam splitter. The pump beam was irradiated

88  on a photoconductive dipole antenna fabricated on a LT-GaAs wafer for generation of the THz

89  waves, and the probe beam was focused onto an electro-optic ZnTe crystal for detection of the

90  THz waves (Taylor. et. al., 2008). The THz pulses emitted by the generator were focused on the

91    sample via two metal parabolic mirrors, and the THz pulses reflected by the sample via two

92    additional parabolic mirrors were guided to the detection antenna. The system measures

93    far-infrared spectra between 0.1 THz and 3.0 THz. The sample was scanned by moving the

94    two-dimensional motorized stage, and the obtained image data were saved and analyzed using a

95    computer. Details about the principles of the system are explained elsewhere (Kim. et. al.,

96    2012). The experiment was performed at room temperature, and the humidity was maintained at

97    approximately zero by purging the system with dry nitrogen to avoid absorption of vapor.

98    *2.2 Sample preparation*

99       Wheat used in the experiment was collected from the School of Food Science and

100   Technology, Henan University of Technology, Zhengzhou, China. The wheat was of the same

101   variety and produced in 2013. Wheat grains were moistened at a humidity of 28% and were

102   evenly distributed in a circular Petri dish. The Petri dish was put into an incubator box that was

103   maintained at a constant temperature of 25°C, where it remained for eight days. Wheat with

104   different stages of mold growth (none, slight, moderate, and serious) where then selected (as

105   shown in Fig.2) and individually imaged by the THz imaging system with a spatial resolution of

106   0.25 mm. For each degree of mold contamination, 50 samples were used without further

107   processing.

108   *2.3 Multivariate Analysis Methods*

109   *2.3.1 Principal component analysis*

110       PCA (Lin, Zhao, Sun, Chen, & Zhou, 2011; Noori, Sabahi, Karbassi, Baghvand, & Zadeh,

111 2010) is a multivariate statistical and dimensional reduction method that can be used to reduce

112 the complexity of input variables when dealing with large datasets. In this method, a large

113 volume of data is transformed into a small number of principal components (PCs). PCs can be

114 expressed as:

$$Z_i = a_{i1}X_1 + a_{i2}X_2 + \cdots + a_{in}X_n \tag{1}$$

116 where $Z_i$ represents the PCs, $a_i$ represents the related eigenvectors, and $X_i$ represents the input

117 variables. This information can be acquired by solving following equation.

$$|R - I\lambda| = 0 \tag{2}$$

119 where $R$ is the variance-covariance matrix, $I$ is the unit matrix, and $\lambda$ is the eigenvector.

120 *2.3.2 Support vector machines*

121 SVM is a widely used, supervised statistical learning method for analyzing data and

122 recognizing patterns (He, Yang, & Xie, 2013; He, Wu, & Sun, 2014). SVM demonstrates

123 advantage over other methods when dealing with small samples, and high-dimensional and

124 non-linear data. In the multi-class SVM method, $k(k-1)/2$ classifiers are constructed, where

125 $k$ is the class number of the data. The following two-class classification problem was

126 implemented by training the *ith* and *jth* data classes:

$$\min_{w^{ij},b^{ij},\xi^{ij}} \frac{1}{2}(w^{ij})^T w^{ij} + c(\sum_t (\xi^{ij})t) \tag{3}$$

128 subject to $\quad\begin{array}{l}(w^{ij})^T \phi(x_t) + b^{ij} \geq 1 - \xi^{ij}, \text{ if } x_t \text{ in } ith \text{ class} \\ (w^{ij})^T \phi(x_t) + b^{ij} \leq -1 + \xi^{ij}, \text{ if } x_t \text{ in } jth \text{ class} \\ \xi^{ij} \geq 0\end{array}$

$$\tag{4}$$

129    where $w$ and $b$ define the optimal hyperplane, $\xi$ represents the slack variable, $c$ is the

130    penalty factor, and $\phi(x)$ is the sample set. Selection of the kernel function in SVM models

131    significantly affects model performance. In this paper, the commonly used radial bias function

132    (RBF) $k(x_i, y_i) = \exp(-\dfrac{\|x_i - y_i\|^2}{\gamma^2})$ was used. The adjustable kernel function parameter $C$ controls

133    the trade-off between the minimum model complexity and minimum training error, while $\gamma$

134    represents the degree of generalization and the width of the kernel function. A grid-search

135    procedure was employed to find the optimal parameters of the model (Maali, & Al-Jumaily,

136    2013).

137        The root mean square error (RMSE) was used to evaluate the performance of the

138    established model (Zhang. et. al., 2008). The RMSE is calculated as

139
$$RMSE = \sqrt{\dfrac{\sum_{i=1}^{N}(y_i^{pre} - y_i)^2}{N}}$$
(5)

140    where $y_i$ represents actual value of the $i$th sample in the data set, $y_i^{pre}$ is the predicted weight

141    ratio value of the $i$th sample in the developed model, and $N$ is the sample size.

142    *2.3.3 Partial least squares regression*

143        PLSR is one of most robust and reliable multivariate-data analysis methods, and is

144    particularly suitable for use in situations where there is a linear relation between the spectra and

145    properties of the considered objects (Brereton, 2000). A PLSR analysis was performed to

146    establish a regression model for the prediction of target chemical concentrations (variable

147    matrix Y) based on the corresponding spectra data (variable matrix X). The underlying PLSR

148    model is expressed as:

149
$$X = TP^T + E$$
$$Y = UQ^T + F$$
(6)

150    where $T$ and $U$ are the feature matrices of the variable matrix of $X$ and $Y$ respectively, $P$

151    and $Q$ represent the orthogonal loading matrices, and $E$ and $F$ are the error terms.

152    *2.3.4 Back propagation neural network*

153    BPNN is a type of nonlinear multi-layer network, and it has been used extensively to solve a

154    variety of classification and regression problems (Dubey, Bhagwat, Shouche, & Sainis, 2006).

155    A BPNN is based on an algorithm that rectifies the weights within each layer in proportion to

156    the error obtained from the previous layer. In this study, an input layer, a hidden layer, and an

157    output layer were used. By optimizing the hidden nodes from the input variables by "trial and

158    error," BPNN was used to classify samples into predefined varieties, and a new output layer

159    that provided a more precise discrimination of a sample's variety was obtained. Details of the

160    BPNN method are discussed extensively elsewhere (Marengo, Bobba, Robotti, & Lenti, 2004).

161    The whole experiment procedure by using THz imaging technique, as illustrated in Fig. 3, is

162    made from three steps to prepare the data structure for mold statuses wheat identification.

163    **3 Results and Discussion**

164    *3.1 Spectral Analysis*

165    *3.1.1 Moldy wheat spectra*

166    After THz images of wheat with different stages of mold growth were acquired, the only

167    wheat grain areas are segmented as the ROIs to exclude the interfering information origin from

168  the background in each image. The spectra of each pixel within the ROI were extracted and

169  averaged at each frequency to generate a mean value, which is then expressed as the ROI

170  spectrum. The average frequency domain spectra of each degree of mold growth, in the range

171  of 0.1–2.0 THz, are shown in Fig. 4. It is seen that an intense trough is present at around 1.67

172  THz, which is related to the absorption of water within the grain. And the spectral curves of

173  these four mold statuses wheat are quite similar at the beginning. Hence, spectral frequency

174  range from 0.2-1.6 THz is employed for further identification study. Meanwhile, the general

175  trends of the four spectral curves show no obvious differences, which indicated that mold

176  statuses of the wheat could not be identified from spectral curves directly.

177  To solve this problem, more sophisticated computational analysis methods were employed

178  to differentiate between the mold statuses of the wheat. Therefore, a dataset with 512 spectral

179  features and 200 wheat samples was selected in order to construct a classification model to

180  discriminate between the different degrees of moldiness. A dataset consisting of 200 samples

181  was randomly split into a calibration set (120 samples) and a prediction set (80 samples). The

182  classification errors would clearly decrease when training more samples. Hence each wheat

183  sample leaves fewer samples to analyze and obtains higher prediction accuracy. But when more

184  training number, redundant information (existed in the large number of input variable) would

185  affect the robust and ability of the classification models. Meanwhile, the less input simplify the

186  classification models and accelerate the calculated speed.

187 *3.1.2 PCA Analysis*

188    PCA was performed on all of the spectral data (with a frequency range of 0.2–1.6 THz)

189 obtained from the normal, slightly moldy, moderately moldy, and seriously moldy wheat

190 samples to reduce the high dimensionality of the problem and qualitatively identify the samples.

191 The explained variance rate for the top four PCs extracted from the original THz spectra data

192 are 93.22%, 3.61%, 1.24%, and 0.21%, respectively. The top four PCs explain 98.25% of the

193 total contribution to the original data. It is shown that the cumulative reliabilities of the top four

194 PCs represent 98% of the total information to the original data. Thus, they contain the

195 maximum information across all the wheat samples and reduce the dimensions from 512

196 spectral measurements for classification of different mold statuses of wheat to only three

197 components. Figure 5 shows the three-dimensional scores plotted for the first three PCs for all

198 of the samples. As we can see, the different mold statuses are distributed separately in the

199 three-dimensional area. However, some sample points near the boundaries of normal and

200 slightly moldy wheat are mixed although their sample points are clustered. Therefore, it is

201 necessary to employ an adequate classification model based on the PCA process for further

202 discrimination.

203 *3.1.3Optimal Frequency Selection*

204    A PCA was used for each ROI image to select the optimal frequencies. PC loadings were

205 employed to identify sensitive frequencies that were highly correlated with each PC. The

206 x-loading weights of the first four PCs were used to select each frequency in the full spectral

207     range. Strong peaks and troughs for the top four PCs were selected as the optimum frequencies.

208     As seen in Fig. 6, six frequencies with the values of 0.32 THz, 0.59 THz, 0.87 THz, 1.0 THz,

209     1.29 THz, and 1.58 THz were selected as discriminators of different moldy statuses. The

210     reduced number of frequencies decreased the time to acquire and process each image.

211     *3.2 Multivariate Data Analysis*

212     *3.2.1 Multivariate Data Analysis Based on Full spectra*

213      SVM, PLSR, and BPNN classification models were used to predict the degree of moldiness

214     using the entire spectral dataset. Within the SVM models, the optimization values for the

215     regularization parameter $\gamma$ and the RBF kernel function parameter $C$ were selected when the

216     smallest RMSE was obtained. The optimal parameters $\gamma$ and $C$ were set at 3.6 and 1.8,

217     respectively, which were determined by using the grid search algorithm. For the BPNN model,

218     after several attempts to optimize the parameters, the learning rate factor, momentum factor,

219     initial weight, permitted training error, and maximal training times were set at 0.1, 0.1, 0.6,

220     0.00001, and 1,000, respectively.

221      The SVM, PLSR, and BPNN models were constructed using the top four PCs as inputs.

222     The discrimination results of normal, slightly moldy, moderately moldy, and seriously moldy

223     wheat in the calibration set and prediction set using these models are presented in Table 1.

224     Table 1 Results of the classification models based on full spectra (Cal. represents the calibration
225                set of the samples and Pre. represents the prediction set of the samples.)

| Model | Accuracy per type (%) | | | | Overall prediction |
|---|---|---|---|---|---|
| | Normal | Slightly moldy | Moderately | Seriously | |

|  | Cal. | Pre. | Cal. | Pre. | moldy Cal. | moldy Pre. | moldy Cal. | moldy Pre. | accuracy (%) |
|---|---|---|---|---|---|---|---|---|---|
| PCA-SVM | 100% | 100% | 100% | 86.67% | 100% | 84% | 100% | 100% | 96.5% |
| PCA-PLSR | 100% | 95% | 91.43% | 86.67% | 88% | 84% | 100% | 95% | 93% |
| PCA-BPNN | 93.33% | 90% | 88.57% | 80% | 84% | 76% | 93.33% | 90% | 87% |

226     As the table shows, the performance of the SVM model was, in general, better than those of

227     the PLSR and BPNN models, and achieved a prediction accuracy of 96.5%. The SVM model

228     achieved a classification rate of the normal and serious moldy statuses of 100% in both the

229     calibration and prediction sets; however, the classification rates of the prediction sets of slightly

230     moldy and seriously moldy wheat were relatively lower. Moreover, the PLSR and BPNN

231     models misclassified some statuses, with an overall prediction accuracy of 93% and 87%,

232     respectively. The results indicate that PLSR and SVM models can be used as effective methods

233     for moldy wheat identification, with the SVM model considered the optimum method.

234     *3.2.2 Multivariate Data Analysis Based on Optimal frequencies*

235     Although the classification models have good moldy wheat prediction performances, the

236     large number of frequency variables resulted in complicated and time-consuming data

237     processing. Instead, the use of optimal-frequency selection can reduce the complexity and time

238     required for model establishment. As a consequence of optimal frequency selection, the top four

239     PCs and the selected six frequencies (0.32 THz, 0.59 THz, 0.87 THz, 1.0 THz, 1.29 THz, and

240     1.58 THz) were used as inputs to the SVM, PLSR, and BPNN models. The performance of the

241    optimized models based only on the optimal frequencies is presented in Table 2.

242    Table 2 Results of the classification models based on their optimal spectra (Cal. represents the

243        calibration set of the samples and Pre. represents the prediction set of the samples.)

| Model | Accuracy per type (%) | | | | | | | | Overall prediction accuracy (%) |
|---|---|---|---|---|---|---|---|---|---|
| | Normal | | Slightly moldy | | Moderately moldy | | Seriously moldy | | |
| | Cal. | Pre. | Cal. | Pre. | Cal. | Pre. | Cal. | Pre. | |
| PCA-SVM | 100% | 100% | 97.14% | 86.67% | 92% | 84% | 100% | 95% | 95% |
| PCA-PLSR | 100% | 95% | 91.43% | 80% | 92% | 84% | 96.67% | 95% | 92.5% |
| PCA-BPNN | 93.33% | 85% | 88.57% | 73.33% | 84% | 76% | 93.33% | 90% | 86% |

244        As shown in Table 2, the BPNN model had the worst prediction result, with a classification

245    accuracy of 86%. The classification rates of the SVM and PLSR models in both the calibration

246    and the prediction sets were all over 80%. The SVM model obtained the highest overall

247    prediction accuracy, 95%, and a classification accuracy of 100% for normal and seriously

248    moldy wheat in the calibration set. The slightly moldy and moderately moldy wheat showed

249    poorer prediction accuracy in all models, compare with the normal wheat and seriously moldy

250    wheat.

251        The plots of the actual values compared to the predicted values using the PCA-SVM

252    models based on the full spectra and selected optimal frequencies are shown Fig. 7. A threshold

253    value (dummy variable ±0.5) was set to define the class limits. Subintervals from 0.5–1.5, 1.5–

254    2.5, 2.5–3.5, and 3.5–4.5 represent normal, slightly moldy, moderately moldy, and seriously

255    moldy wheat samples, respectively. It can be seen in Figs. 6 (a) and (b) that a similar

256    distribution of points between the full spectrum and the optimal frequencies was obtained. The

257    experimental results demonstrate the feasibility of using selected optimal frequencies for the

258    discrimination of wheat grains with different mold statuses.

259    *3.3 THz Images of Moldy Wheat*

260        The implementation of a visualization process is helpful for determining the degree of

261    moldiness of a wheat grain, which can be difficult when observed by just the naked eye. In this

262    study, the PCA-SVM model acquired the best classification accuracy and therefore was used to

263    generate THz moldy wheat images. Training of the SVM model was done using the optimal

264    frequencies selected by the PCA. The reduced spectral data were then used as input to the SVM

265    model. The output value of the model was the reflectivity of each pixel, which corresponds to a

266    different component within each wheat grain. When the values of all pixels within the wheat

267    grain were calculated, an image was generated based on the spatial positions of each pixel.

268        Figure 8 shows the THz images of normal, slightly moldy, moderately moldy, and seriously

269    moldy wheat. Regions (1), (2), and (3) represent the embryo of each wheat grain. Except for the

270    embryo structure, the inner structures of the wheat sample in Fig. 7(a) and 7(b) are evenly

271    distributed. However, in Fig. 7(b) the embryo and edge structure have changed, indicating that

272    the wheat is in its moldy infancy, while it is seen that the wheat in Fig. 7(a) is not contaminated

273    with mold. In Fig. 7(c), the embryo area and small range of inner structures are damaged,

274    indicating that the sample has a moderate degree of mold growth. Finally, in 7(d), the red area

275　(5) indicates that the inner structures of this wheat sample are totally damaged, and the

276　embryonic area is absent.

277　*3.4 Discussion*

278　　The excellent discrimination results demonstrate that the THz reflection imaging technique

279　combined with PCA feature extraction and a SVM classification model can be used to identify

280　wheat grains with different mold statuses. Six optical frequencies (0.32 THz, 0.59 THz, 0.87

281　THz, 1.0 THz, 1.29 THz, and 1.58 THz) were selected according to the top four PC loading

282　weights. The overall prediction accuracy of the PCA-SVM model based on the selected optimal

283　frequencies was 95%, which is higher than that achieved with the PCA-PLS and PCA-BPNN

284　models. The optimal frequency-based models used six frequencies instead of 159 frequencies,

285　indicating a decrease of 96.49%. The performance of each classification model showed only a

286　slight decline from full spectra to optimal frequencies, implying that the optimal frequencies

287　were effective, and as such, we encourage further study of them. Furthermore, the fewer input

288　variables accelerated the data calculation speed and simplified the model complexity. In further

289　studies, different frequency selection methods and different classification models will be

290　applied to improve the prediction accuracy and explore the optimal frequency for moldy wheat

291　identification.

292　　Additionally, the PCA-SVM model was used to classify the THz image data and determine

293　the degree of mold contamination as normal, slightly moldy, moderately moldy, and seriously

294　moldy. The THz images provided information regarding the spatial distribution of different

295     components within the wheat grain, and were helpful for detecting changes in a grain's inner

296     structure due to varying mold status. Our results show that THz imaging can be used to

297     recognize the wheat when it is in its early moldy stage, which cannot be done with conventional

298     imaging and spectroscopy, and thus provides an early warning technique for mold

299     contamination. The THz imaging technique has the potential to be an effective tool for

300     agriculture quality and safety control. Therefore, it is essential to expand the sample variety

301     number and optimize the image classification algorithm in further studies to assist

302     in discriminating the multiple statuses of wheat mold en masse and for practical applications.

303     **4 Conclusion**

304     THz imaging combined with multivariate data analyses was employed to discriminate

305     wheat grains with different mold statuses. Spectral information was extracted from the THz

306     images, in the range of 0.2–1.6 THz, for each wheat sample. The feature data of each spectrum

307     were analyzed and six optimal frequencies were selected using PCA. In addition, the SVM,

308     PLSR, and BPNN models were constructed based on the full spectra and optimal frequencies to

309     help discriminate between different moldy wheat samples. The prediction accuracies of the full

310     spectra were similar to those obtained using only the optimal frequencies. The PCA-SVM

311     model was considered to be the optimal model, and the prediction accuracies reached 95%. The

312     PCA-SVM model was also used on THz images as a visual demonstration of the classification

313     technique. Our experimental results demonstrate that THz imaging is a potential tool for the

314     classification of moldy wheat.

**RSC Advances Accepted Manuscript**

324   **REFERENCES**

325   Arngren, M., Hansen, P. W., Eriksen, B., Larsen, J., & Larsen, R. (2011). Analysis of

326   Pregerminated Barley Using Hyperspectral Image Analysis. *Journal of Agricultural and Food*

327   *Chemistry*, 59(21), 11385-11394.

328   Brereton, R. G. (2000). Introduction to multivariate calibration in analytical chemistry. *Analyst*,

329   125(11), 2125-2154.

330   Dubey, B. P., Bhagwat, S. G., Shouche, S. P., & Sainis, J. K. (2006). Potential of artificial neural

331   networks in varietal identification using morphometry of wheat grains. *Biosystems Engineering,*

332   95(1), 61-67.

333   Eifler, J., Martinelli, E., Santonico, M., Capuano, R., Schild, D., & Di Natale, C. (2011).

334   Differential Detection of Potentially Hazardous Fusarium Species in Wheat Grains by an

335   Electronic Nose. *Plos One*, 6(6).

336   EIMasry, G., Wang, N., EISayed, A., Ngadi, M., 2007. Hyperspectral imaging for nondestructive

337   determination of some quality attributes for strawberry. Journal of Food Engineering, 81 (1), 98–

338   107.

339   Ferguson, B., & Zhang, X. C. (2002). Materials for terahertz science and technology. *Nature*

340   *Materials,* 1(1), 26-33.

341   Fukunaga, K., & Hosako, I. (2010). Innovative non-invasive analysis techniques for cultural

342   heritage using terahertz technology. *Comptes Rendus Physique*, 11(7-8), 519-526.

343   Ge, H. Y., Jiang, Y. Y., Xu, Z. H., Lian, F. Y., Zhang, Y., & Xia, S. H. (2014). Identification of

344    wheat quality using THz spectrum. *Optics Express*, 22(10), 12533-12544.

345    Gowen, A. A., O'Sullivan, C., & O'Donnell, C. P. (2012). Terahertz time domain spectroscopy

346    and imaging: Emerging techniques for food process monitoring and quality control. *Trends in*

347    *Food Science & Technology,* 25(1), 40-46.

348    Guillet, J. P., Recur, B., Frederique, L., Bousquet, B., Canioni, L., Manek-Honninger, I.,

349    Desbarats, P., & Mounaix, P. (2014). Review of Terahertz Tomography Techniques. *Journal of*

350    *Infrared Millimeter and Terahertz Waves*, 35(4), 382-411.

351    He, M., Yang, G. L., & Xie, H. Y. (2013). A hybrid method to recognize 3D object. *Optics*

352    *Express*, 21(5), 6346-6352.

353    He, H. J., Wu, D., & Sun, D. W. (2014). Potential of hyperspectral imaging combined with chem

354    ometric analysis for assessing and visualising tenderness distribution in raw farmed salmon fillet

355    s. Journal of Food Engineering, 126, 156-164.

356    Kim, K. W., Kim, K. S., Kim, H., Lee, S. H., Park, J. H., Han, J. H., Seok, S. H., Park, J., Choi,

357    Y., Kim, Y. I., Han, J. K., & Son, J. H. (2012). Terahertz dynamic imaging of skin drug

358    absorption. *Optics Express,* 20(9), 9476-9484.

359    Lin, H., Zhao, J. W., Sun, L., Chen, Q. S., & Zhou, F. (2011). Freshness measurement of eggs

360    using near infrared (NIR) spectroscopy and multivariate data analysis. *Innovative Food Science*

361    *& Emerging Technologies*, 12(2), 182-186.

362    Maali, Y., & Al-Jumaily, A. (2013). Self-advising support vector machine. *Knowledge-Based*

363    *Systems*, 52, 214-222.

364 Marengo, E., Bobba, M., Robotti, E., & Lenti, M. (2004). Hydroxyl and acid number prediction

365 in polyester resins by near infrared spectroscopy and artificial neural networks. *Analytica*

366 *Chimica Acta*, 511(2), 313-322.

367 Melinger, J. S., Laman, N., & Grischkowsky, D. (2008). The underlying terahertz vibrational

368 spectrum of explosives solids. *Applied Physics Letters,* 93(1).

369 M.Tonouchi. (2007). Cutting-edge terahertz technology. *Nature Photonics*, 1(2), 97-105.

370 Neethirajan, S., Karunakaran, C., Jayas, D. S., & White, N. D. G. (2007). Detection techniques

371 for stored-product insects in grain. *Food Control*, 18(2), 157-162.

372 Noori, R., Sabahi, M. S., Karbassi, A. R., Baghvand, A., & Zadeh, H. T. (2010). Multivariate

373 statistical analysis of surface water quality based on correlations and variations in the data set.

374 *Desalination*, 260(1-3), 129-136.

375 Oh, S. J., Kim, S. H., Ji, Y. B., Jeong, K., Park, Y., Yang, J., Park, D. W., Noh, S. K., Kang, S.

376 G., Huh, Y. M., Son, J. H., & Suh, J. S. (2014). Study of freshly excised brain tissues using

377 terahertz imaging. Biomedical *Optics Express*, 5(8), 2837-2842.

378 Oladunmoye, O. O., Akinoso, R., & Olapade, A. A. (2010). Evaluation of Some

379 Physical-Chemical Properties of Wheat, Cassava, Maize and Cowpea Flours for Bread Making.

380 *Journal of Food Quality*, 33(6), 693-708.

381 Reid, C. B., Pickwell-MacPherson, E., Laufer, J. G., Gibson, A. P., Hebden, J. C., & Wallace, V.

382 P. (2010). Accuracy and resolution of THz reflection spectroscopy for medical imaging. *Physics*

383 *in Medicine and Biology*, 55(16), 4825-4838.

384    Siegel, P. H. (2004). Terahertz technology in biology and medicine. *Ieee Transactions on*

385    *Microwave Theory and Techniques,* 52(10), 2438-2447.

386    Safrai, E., Ben Ishai, P., Polsman, A., Einav, S., & Feldman, Y. (2014). The Correlation of ECG

387    Parameters to the Sub-THz Reflection Coefficient of Human Skin. *Ieee Transactions on*

388    *Terahertz Science and Technology*, 4(5), 624-630.

389    Taylor, Z. D., Singh, R. S., Culjat, M. O., Suen, J. Y., Grundfest, W. S., Lee, H., & Brown, E. R.

390    (2008). Reflective terahertz imaging of porcine skin burns. *Optics Letters*, 33(11), 1258-1260.

391    Turner, N. W., Subrahmanyam, S., & Piletsky, S. A. (2009). Analytical methods for

392    determination of mycotoxins: A review. *Analytica Chimica Acta*, 632(2), 168-180.

393    Wang, S. J., Liu, K. S. , Yu, X. J. , Wu, D., & He, Y.   2011. Application of hybrid image

394    features for fast and non-invasive classification of raisin. Journal of Food Engineering, 109 (3),

395    531-537.

396    Zhang, Y., Peng, X. H., Chen, Y., Chen, J., Curioni, A., Andreoni, W., Nayak, S. K., & Zhang,

397    X. C. (2008). A first principle study of terahertz (THz) spectra of acephate. *Chemical Physics*

398    *Letters,* 452(1-3), 59-66.

399

400

**Figure captions**

Fig.1 THz reflectance imaging experimental setup.

Fig.2 Wheat samples with different stages of mold contamination: (a) normal; (b) slightly; (c)

moderately; (d) seriously.

Fig. 3 Flowchart of the procedure of discrimination moldy wheat by using THz imaging: (a)

Imaging pre-processing; (b) Spectral analysis; (c) Imaging visualization.

Fig. 4 Frequency-domain THz spectra of the moldy wheat samples

Fig.5 Scores scatter plot of PC1, PC2, and PC3 for each moldy wheat sample

Fig.6 Loading weights of the top four PCs used for selecting the optimal frequencies

Fig.7 Scatter plots of the actual value versus the predicted value using the PCA-SVM model

based on (a) the full spectrum and (b) the optimal frequencies for different moldy wheat samples.

Fig. 8 THz images of four wheat grains with different mold statuses: (a) normal; (b) slightly

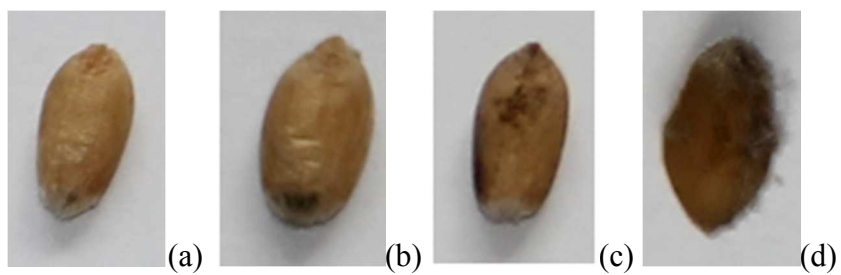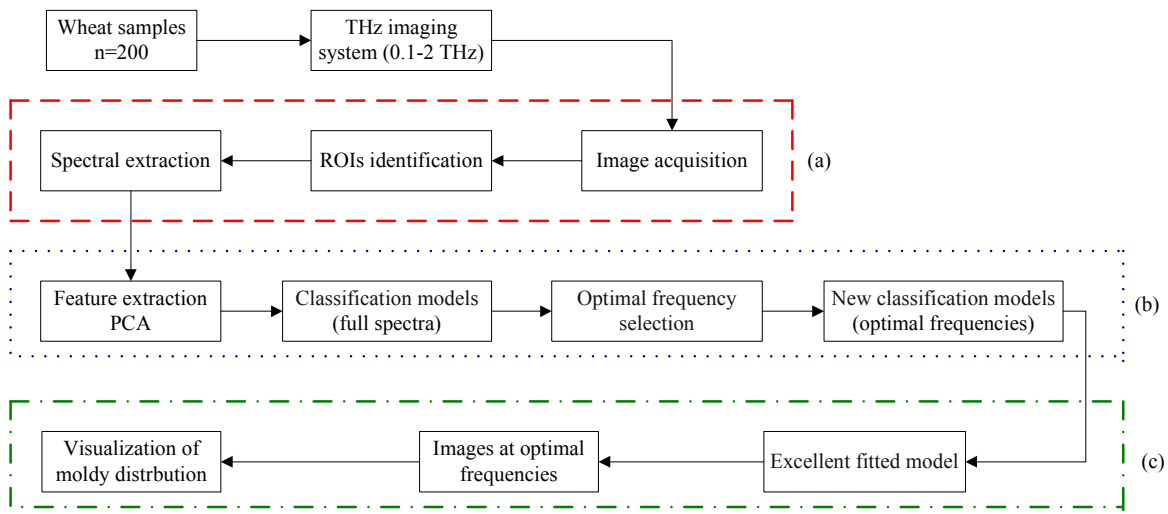moldy; (c) moderately moldy; (d) seriously moldy.
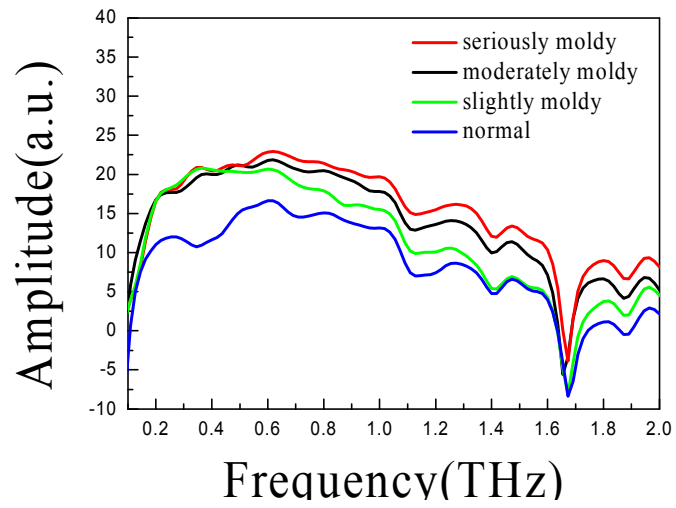
# Figures

Fig. 1

Fig. 2



(a)       (b)       (c)       (d)

Fig. 3

Fig. 4

Fig. 5

Fig. 6

Fig. 7

Fig. 8



(a)    (b)    (c)    (d)