

# Machine Learning Analysis of Phage Oxidation for Rapid Verification of Wash Water Sanitation

Hemiao Cui <sup>a</sup>, Reza Ovissipour<sup>a†</sup>, Xu Yang<sup>a‡</sup>, Nitin Nitin<sup>a,b,\*</sup>

4<sup>a</sup> Department of Food Science and Technology, University of California-Davis, Davis, CA  
595616, United States

6<sup>b</sup> Department of Agricultural and Biological Engineering, University of California-Davis, Davis,  
7CA 95616, United States

13Running Title: Machine learning analysis of phage oxidation for verification of wash water  
14sanitation

15\* Address correspondence to Nitin Nitin, [nnitin@ucdavis.edu](mailto:nnitin@ucdavis.edu)

16<sup>†</sup>Present address: Reza Ovissipour, Seafood AREC/Department of Food Science and  
17Technology, Virginia Polytechnic Institute and State University, Hampton, VA 23669, USA

18‡Current address: Nutrition and Food Science Department, California State Polytechnic  
19University Pomona, Pomona, CA, United States of America

20'Declarations of interest: none

23

24

## 25Abstract

26 The current approaches for process verification during sanitation of fresh produce and other  
27minimally processed products are limited to point measurements of sanitizer concentration at  
28discrete locations and lack rapid biological measurements to assess effectiveness of sanitation.  
29To address this gap, this study evaluates immobilized T7 phage on anodisc membrane  
30(phage@anodisc) as a surrogate for process verification. Fourier Transform infrared (FTIR)  
31spectroscopy results suggested that both chlorine and Peracetic acid (PAA) caused phage DNA  
32damage and protein oxidation. The Gradient Boosting algorithm was employed to develop  
33predictive model for sanitizer concentration levels and *Escherichia coli* O157:H7 inactivation.  
34The machine learning model predicted both the effective sanitizer concentration level and  
35bacterial reduction with ROC (receiver operating characteristic) values between 0.86 and 0.93.  
36Overall, this study identified spectral measurement of phage particles in combination with  
37machine learning approach as an effective tool for process verification.

38

39 **Keywords:** Fresh produce, sanitation, rapid verification, phage, vibrational spectroscopy,  
40machine learning

41

42

## 431. Introduction

3

2

4

44 Fresh produce safety is recognized as one of the challenges by the U.S. Food and Drug  
45Administration. Food contamination by microorganisms may occur at various stages in the food  
46supply chain. Postharvest handling of fresh produce usually involves various cooling and  
47washing steps as well as various mechanical equipment for transportation, storage and packaging  
48of fresh produce. During these handling steps, fresh produce can be contaminated with microbes  
49from wash water or food contact surfaces (Suslow, 1997). Hence, disinfection of wash water and  
50equipment is a critical step to ensure the safety and quality of fresh produce (Suslow, 1997; Gil,  
51Selma, Lopez-Galvez, & Allende, 2009; Cossu, Le, Young, & Nitin, 2017). Monitoring and  
52rapid validation of sanitization efficacy is critical for fresh produce industry to provide safe  
53products, and meet the preventive control requirements of the Food Safety Modernization Act  
54(Brackett, Ocasio, Waters, Barach, & Wan, 2014). The current approaches for validation of  
55sanitation include the standard plate counting methods, water chemistry based on sanitizer  
56concentration, total organic content, oxidation reduction potential (ORP), turbidity and pH of the  
57aqueous phase (Cossu et al., 2017). However, these methods are limited in direct assessment of  
58biological damage induced by sanitizers and can be influenced by complexity of the  
59environment, such as fouling of electrodes and the presence of organic matter (Cossu et al.,  
602017).

61 To assess the biological response to sanitizers in wash water, previous studies have explored  
62measurement of changes in cell membrane permeability, enzymatic activity, protein oxidation  
63and DNA damage (Cossu et al., 2017). These measurements suggest the potential for measuring  
64biochemical changes to validate sanitation efficacy. One of the key limitations of this approach is  
65that the biochemical measurements required multiple sample preparation and biochemical  
66reactions steps to assess changes induced by sanitizers. Furthermore, due to the constraints of

67introducing live bacteria in food facilities including commensal bacteria, translation of these  
68biological approaches for measurements in food industry is limited. In contrast to these  
69biochemical assays, we recently developed a spectroscopic approach to quantify oxidation of  
70isolated DNA molecules after exposing to different chlorine concentrations using vibrational  
71spectroscopy and chemometrics (Ovissipour, Rai, & Nitin, 2019). Spectroscopic approach  
72significantly reduces the operational complexity of biochemical measurements in cells both by  
73reducing the time and the number of manual steps required for the biochemical assay. Thus,  
74spectroscopic analysis combined with appropriate biological surrogate can provide a platform to  
75rapidly assess sanitation efficacy. The spectroscopic analysis using machine learning approaches  
76can reduce the multidimensional spectroscopic dataset and uncover complex relationship (Zareef  
77et al., 2020). Supervised machine learning models can be used for quantitative analysis to  
78identify biochemical changes in biological surrogates. Predictive quantitative models can be built  
79using proper reference data set such as sanitizers concentrations and bacterial population.

80 The overall objectives of this study were to evaluate the potential of using phage as a  
81biological surrogate for assessment of sanitation efficacy of wash water using vibrational  
82spectroscopic measurements. Bacteriophage was selected as a biological surrogate as phages are  
83commonly present in the environment, relatively easy amplification procedures to generate  
84phages and simple structural composition (nucleic acid and protein). In addition, phages are  
85commonly used as an indicator organism to evaluate contamination of water and likely to be  
86more widely acceptable by food industry as a surrogate compared to live bacteria. Upon  
87interaction between phages and sanitizers (e.g. chlorine or peracetic acid), the results of phage  
88DNA oxidation and DNA conformational changes induced by wash water with chlorine or  
89peracetic acid were measured using FTIR. Spectra from FTIR were analyzed using both principle

component analysis for classification of the spectral data and the Gradient Boosting Algorithm for quantitative predictive model development. The results of this study illustrate the potential of a novel approach to validate antimicrobial potential of wash water in the fresh produce industry.

## 2. Materials and methods

### Reagents and supplies

Peroxyacetic acid (PAA), NaOCl (10% sodium hypochlorite), citric acid, Dey-Engley neutralizing broth (D-E broth), tryptic soy broth (TSB) and tryptic soy agar (TSA) were obtained from Sigma-Aldrich (St. Louis, MO, USA). Whatman® anodisc inorganic filter membrane (13 mm, 0.02 µm pore size) was obtained from GE Healthcare (Buckinghamshire, UK). Filtration system (a filtering flask and a fritted glass base) was obtained from Fisher Scientific (Pittsburgh, PA, USA). Phosphate buffered saline (PBS) was purchased from Fisher Bioreagents (Fair Lawn, NJ, USA). Milli-Q water was produced by QPAK® 2 purification system (EMD Millipore, Billerica, MA, USA).

### Bacterial cultures and phage preparation

Both *E. coli* BL21 (ATCC BAA-1025) and bacteriophage T7 (ATCC BAA-1025-B2) were obtained from the American type culture collection. A shiga-toxin knockout rifampicin-resistant *E. coli* O157:H7 mutant (ATCC 700728) was kindly given by Dr. Linda Harris (University of California, Davis). Both *E. coli* strains were cultured in TSB broth at 37 °C for 16 hours before use.

Bacteriophages were propagated as the following procedure. Bacteriophages were first inoculated into log-phase *E. coli* BL21 culture at the ratio of 1:100 (phage: bacteria). The mixture was incubated at 37 °C for 15 min for initial infection and then centrifuged at  $16100 \times g$  for 10 min. Supernatant was discarded and the same volume of TSB was added to resuspend the

113pellet, followed by incubating at 37 °C with 200 rpm shaking until no visible turbidity was  
114observed. Chloroform was then added to the final concentration at 20% (vol/vol) and incubated  
115at 4 °C overnight. Then, chloroform added mixture was centrifuged at  $5,000 \times g$  for 10 min and  
116water phase was collected. The water phase which contained free phages at around  $10^{10}$  PFU/ml  
117was used for further filtration and FTIR analysis.

#### 118**Preparation of sanitizer solutions**

119Both PAA and NaOCl and their concentrations were selected based on their application in food  
120industry. PAA was prepared at final concentrations of 20, 40, 60 or 80 ppm while NaOCl was  
121prepared at final free NaOCl concentrations of 2, 5, or 10, 15 ppm and adjusted pH to 6.5-7.0  
122using 0.5 M citric acid.

#### 123**Preparation of phage@anodisc**

124Free phage solution (1 ml) prepared as previously described was gently pipetted onto anodisc  
125and filtered using the filtration system. After filtration, 1 ml of sterile MilliQ water was gently  
126added to anodisc to wash and get rid of soluble impurities that generated during phage  
127preparation.

#### 128**Exposure of phage@anodisc to PAA or NaOCl**

129Five phage@anodisc were exposed to PAA at 0, 20, 40, 60 or 80 ppm for 2 min at 4 °C. The time  
130and temperature were selected based on food industry sanitation protocol. Similarly, another five  
131phage@anodisc were also exposed to 0, 2, 5, 10, 15 ppm of NaOCl solution at the same  
132condition. All phage@anodisc were then rinsed with 0.1% sodium thiosulfate to inactivate the  
133NaOCl followed by MilliQ water, before FTIR analysis.

#### 134**Inactivation of bacteriophage or *E. coli* using PAA or NaOCl**

135Bacteriophages or *E. coli* O157:H7 cells were inoculated into PAA or NaOCl solutions with  
136different sanitizer concentration levels following the same procedure as described earlier for 2  
137min at 4 °C. Survivor population of both phage T7 and *E. coli* O157:H7 were enumerated to  
138quantify inactivation of these microbes as a function of sanitizer concentration.

139Briefly, *E. coli* O157:H7 cells treated with PAA were 10-fold serially diluted in D-E broth and  
140allow for incubating at room temperature for at least 10 min before plating to allow recovery of  
141injured cells. In comparison, *E. coli* O157:H7 cells treated with NaOCl were first neutralized  
142using 0.1 M sodium thiosulfate before serial dilution and plating. Enumeration of phage T7  
143particles was conducted by co-incubation of 10-fold diluted phage suspension and its host *E. coli*  
144BL21 in soft TSA agar (0.75%), followed by gently pouring the mixture in empty petri-dishes.  
145The clear plaque can be observed and enumerated after incubation at room temperature  
146overnight.

#### 147**Fourier Transform infrared spectroscopy (FTIR)**

148Phage@anodisc exposed to selected levels of PAA and NaOCl concentrations were dried under  
149the laminar hood for 2 h. FTIR spectra were collected from 4000 to 400 cm<sup>-1</sup> at a resolution of 2  
150cm<sup>-1</sup> from the phage@ anodisc membrane samples (32 interferograms) (IRPrestige-21 FTIR  
151spectrometer, Shimadzu Co., Kyoto, Japan).

#### 152**FTIR data modelling**

153The FTIR data was pre-processed with Python programming language [to normalize and](#)  
154[concatenate FTIR measurements](#). Principal component analysis PCA was then conducted with  
155the scikit-learn library PCA package using the FTIR data as an input [to reduce dimensionality of](#)  
156[the data set and identify the principal components for classification of the spectroscopy data set](#).  
157(Li & Phung, 2014).

158      The relationships between spectral signals and sanitizers concentrations as well as  
159bacterial inactivation were also modeled using the Gradient Boosting algorithm, specifically the  
160LightGBM. LightGBM is a Gradient Boosting Decision Tree (GBDT) based machine learning  
161algorithm. Given a set of training data  $(x_1, y_1), \dots, (x_n, y_n)$  where  $y$  refers to sanitizer dosage/bacteria  
162log reduction and  $x$  refers to spectral signal at the same treatment condition, the optimization  
163goal of GBDT is to minimize the loss function which quantifies the prediction error. Details of  
164GBDT is shown below:

Algorithm detail for Boosting Regression Trees
--

1. Set $\widehat{f}(x) = 0$ and $r_{(i)} = y_{(i)}$ for all $i$ in the training set where $r_{(i)}$ represents <i>ith</i> residual
--

2. For each tree $b = 1, 2, \dots, B$ where $B$ is the total number of trees, repeat:
---

a. Fit a tree $\widehat{f}^b$ with $d$ splits to the training data $(X, r)$
---

b. Update $\widehat{f}$ by adding in a shrunk version of the new tree where $\lambda$ is the learning rate:
---

$\widehat{f}(x) \leftarrow \widehat{f}(x) + \lambda \widehat{f}^b$
--

c. Update residuals,
----------------------

$r_{(i)} \leftarrow r_{(i)} - \lambda \widehat{f}^b$
--

3. Output the boosted model,
------------------------------



$$\widehat{f(x)} = \sum_b \lambda \widehat{f^b}$$

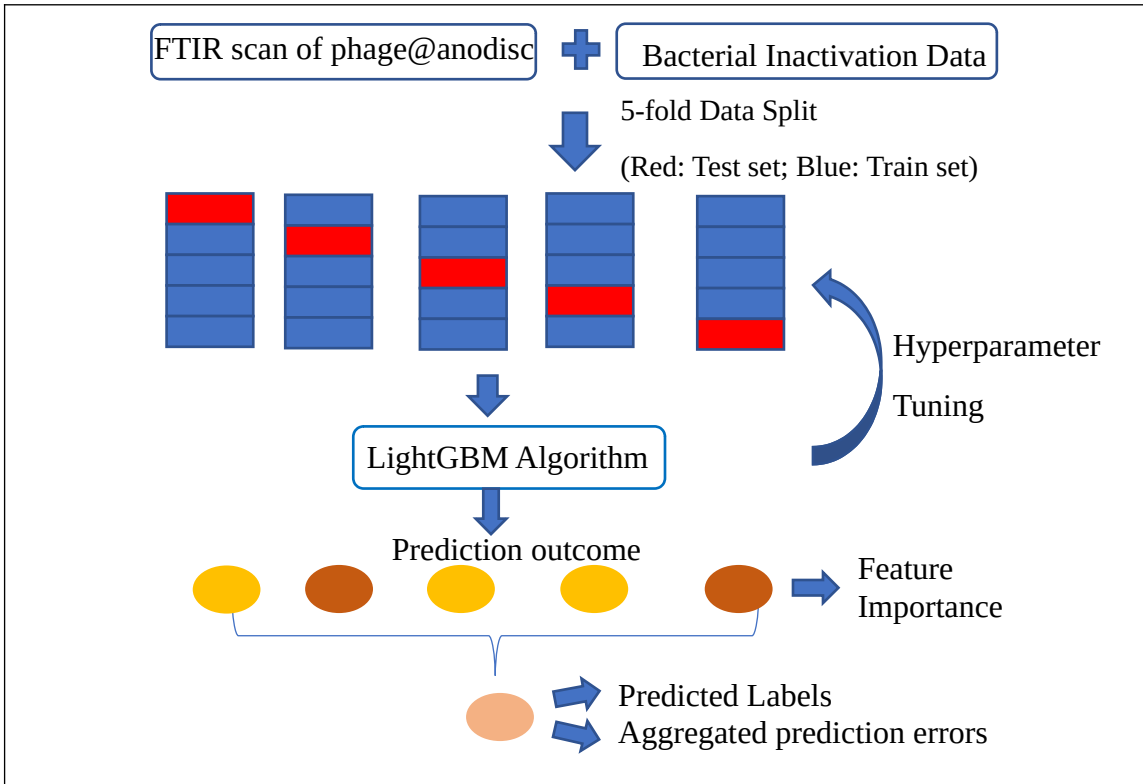
165

166LightGBM builds on top of GBDT by applying Gradient-based One-Side Sampling (GOSS)  
 167method and exclusive feature bundling (EFB). LightGBM uses GOSS to determine the split  
 168point via calculating variance gain, thereby accelerating the training for each level of the tree. It  
 169uses EFB to bundle exclusive features into a single feature to reduce the number of features and  
 170uses greedy algorithm for approximation of the best split point (Ke et al., 2017). These render  
 171LightGBM with ~~was chosen due to its~~ fast speed and the ability to ~~handle a~~ high dimensional  
 172large dataset. FTIR data tends to be high dimensional. Thus, LightGBM algorithm was selected  
 173to predict sanitizers concentrations as well as bacterial inactivation from FTIR spectra from.

174~~Specifically, it uses the Gradient-based One-Side Sampling method to exclude a significant~~  
 175~~proportion of data instances with small gradients, thereby accelerating the training for each level~~  
 176~~of the tree. In addition, LightGBM bundles mutually exclusive features (rarely take nonzero~~  
 177~~values simultaneously) to reduce the number of features and uses greedy algorithm for~~  
 178~~approximation (Ke et al., 2017).~~

179LightGBM model was implemented with the scikit-learn library. **Figure 1** showed the prediction  
 180pipeline. FTIR data (predictor variables) and bacterial inactivation data (response variables) were  
 181concatenated. To investigate the model performance to unseen data set, 5-fold cross-validation  
 182was conducted by splitting the combined data set to 5 folds. For each fold, model was trained on  
 183training datasets (represented as blue cells in Figure 1) and evaluated on test datasets  
 184(represented as red cells in Figure 1). The default loss function (multi Log loss) was used to  
 185construct the objective function. Predictions was evaluated based on the Receiver Operating

186Characteristic (ROC) curve and Confusion Matrix (CM). Prediction errors for all 5 folds were  
187aggregated to give the final ROC and CM value in the figures.- [Sample Python code for the](#)  
188[analysis was provided.](#)



189

190

**Figure 1** Flowchart of the model development pipeline

191

**192Statistical analysis**

193Throughout the study, experiments were conducted in three independent trials. For each trial,  
194measurements were conducted in three replicates. The significant differences between treatments  
195were determined through one-way Analysis of Variance (ANOVA) followed by Tukey's  
196pairwise comparisons and  $p < 0.05$  is considered as significant.

19

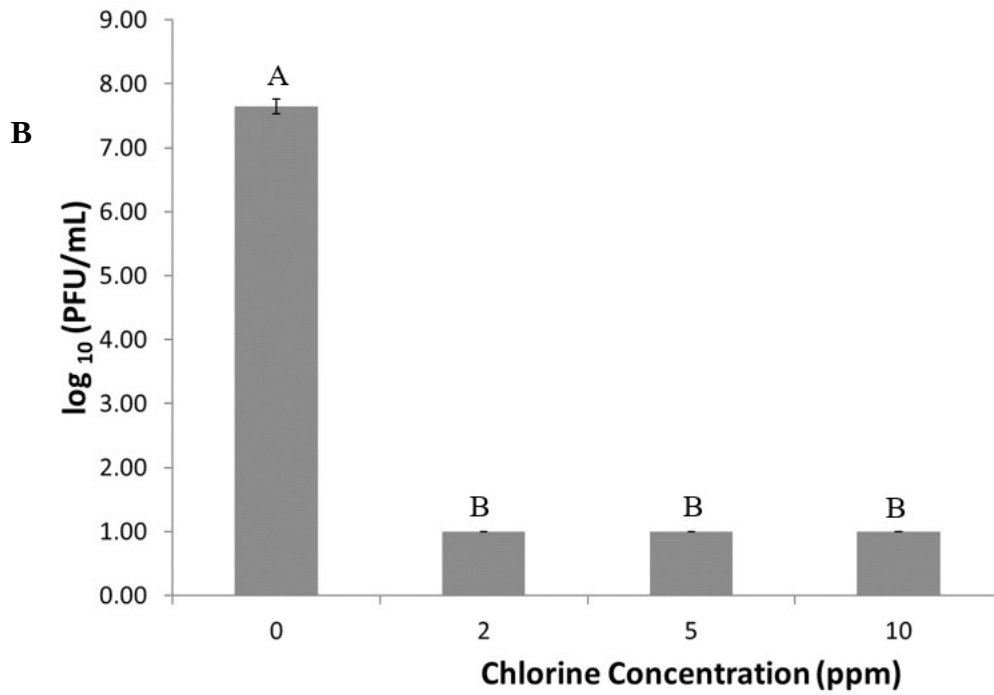
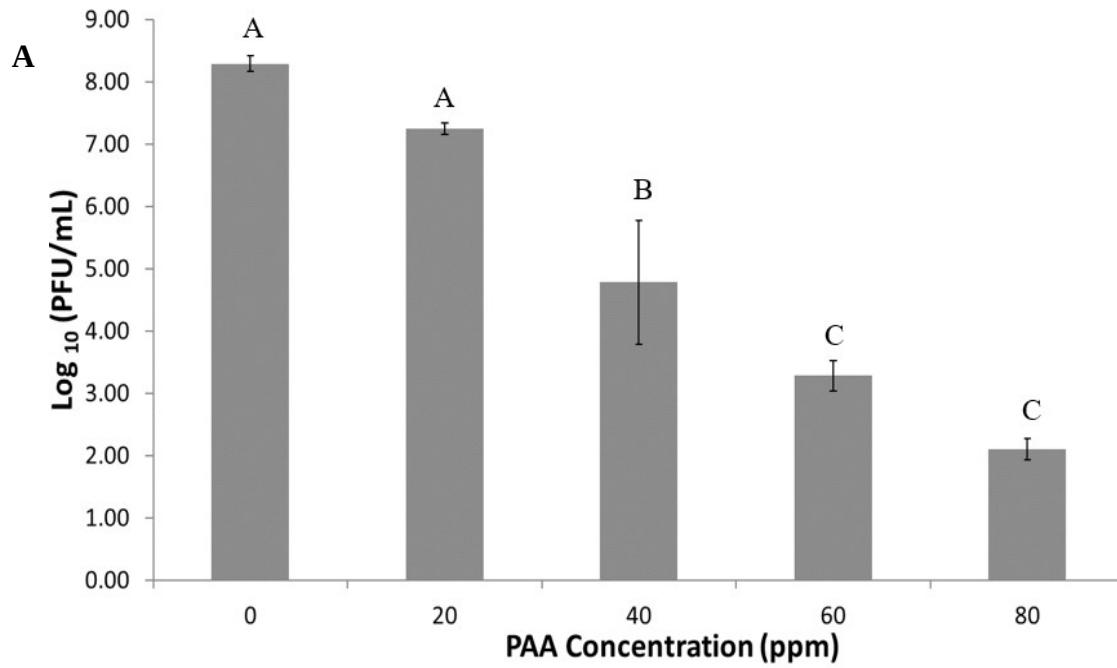
20

### 1973. Result and discussion:

#### 198Inactivation of phage T7 by selected sanitizers

199**Figure 2** illustrates the survivor population of phage T7 upon treatment with PAA or NaOCl at  
200varying levels of sanitizer concentration. As shown in **Fig 2A**, inactivation of phage T7 increased  
201with an increase of PAA concentration ( $P < 0.05$ ). PAA at 80 ppm concentration successfully  
202inactivated more than 6-log of phage T7 but despite using high concentration levels, complete  
203inactivation of phages (9 log inoculum level) was not observed as shown in **Fig 2A**. In  
204comparison, complete inactivation of inoculated phages was observed at even the lowest  
205concentration tested (2 ppm of free NaOCl) in this study. The results suggested that viral  
206particles such as bacteriophages may be more susceptible to NaOCl than PAA. Similar results  
207have also been reported by Morin *et al.* that MS2 phages are more resistant to PAA than NaOCl  
208(Morin et al., 2015). This difference may result due to broad reactivity of NaOCl with proteins  
209and DNA molecules, thus damaging both the capsid proteins and the DNA of viral particles.

210



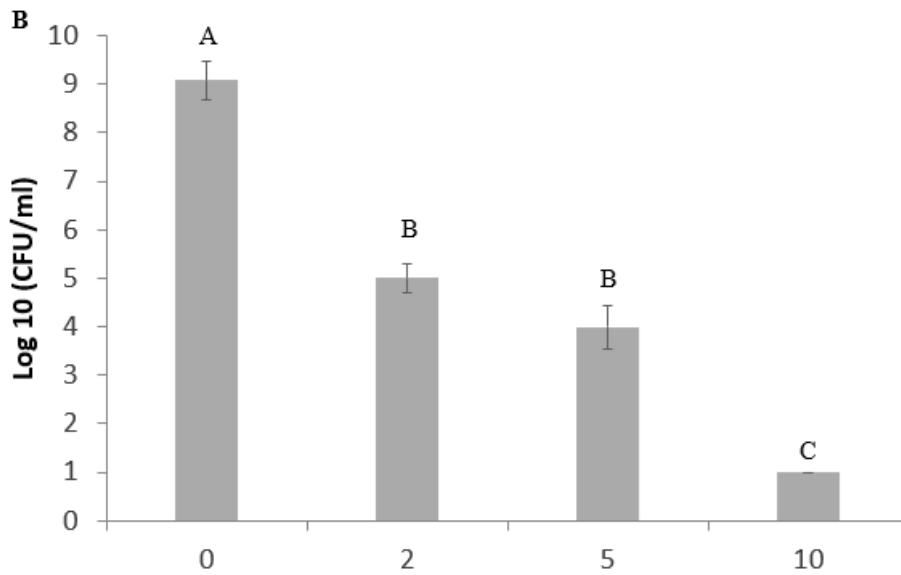
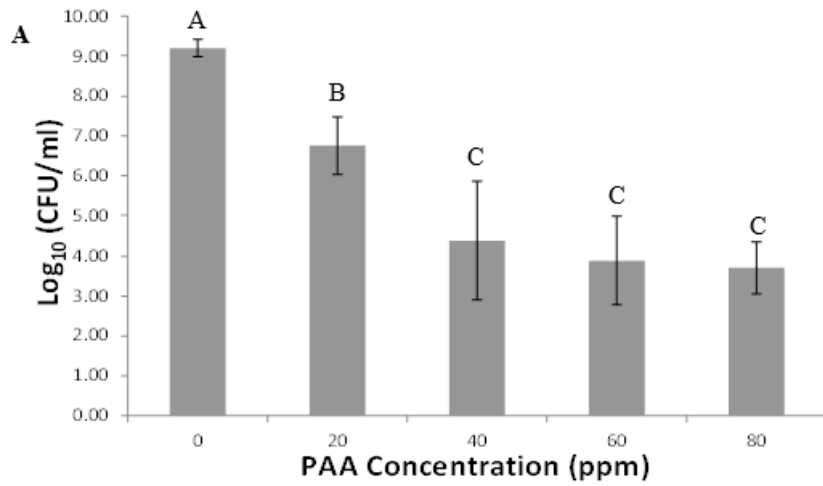
215 **Figure 2** Survivor plot of phage treated by PAA (A) or NaOCl (B) at selected  
216 concentrations for 2 min at 4 °C.

217

#### 218 **Inactivation of *E. coli* O157:H7 by selected sanitizers**

219 **Figure 3** illustrates the survivor population of *E. coli* O157:H7 upon treatment with PAA or  
220 NaOCl at varying levels of sanitizer concentration. As shown in **Fig 3A**, *E. coli* O157:H7 cells  
221 were significantly inactivated (2-log inactivation) by PAA, even at 20 ppm. However, no  
222 significant increase in inactivation of *E. coli* O157:H7 was observed with an increase of PAA  
223 concentration above 40 ppm. Even at the highest levels of PAA (80 ppm for 2 min) used in this  
224 study, only 5 log inoculated bacterial cells were inactivated from the initial inoculum levels of 9  
225 log of bacteria. **Fig 3B** showed *E. coli* O157:H7 reduction upon treatment with different  
226 concentration levels of NaOCl after 2 min at 4°C. NaOCl at the levels of 2 and 5 ppm caused 4  
227 and 5 log reduction, and 10 ppm of NaOCl completely inactivated *E. coli*.

228



229

230

231

232 **Figure 3** Survivor plot of *E. coli* O157:H7 treated by PAA (A) or NaOCl (B) at selected  
 233 concentrations for 2 min at 4 °C.

234

235 **PCA models**

27

14

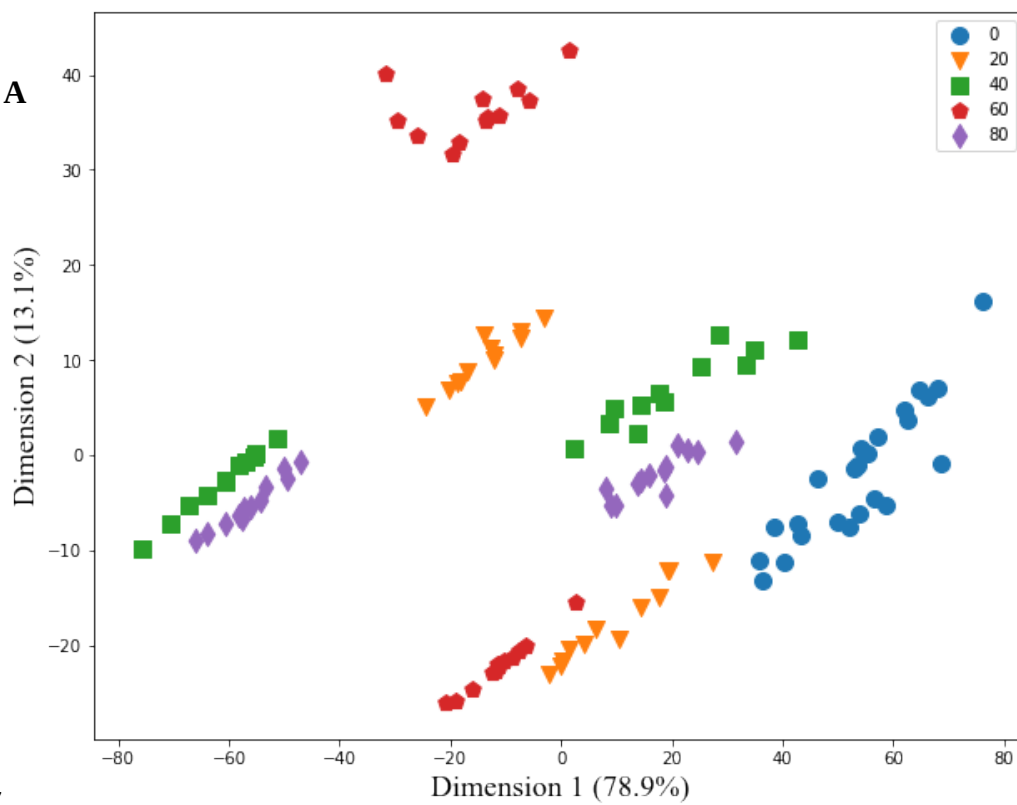
28

236PCA analysis has been used to describe the variations in evaluating oxidation of DNA in live *E.*  
237*coli* cells upon exposure to sanitizers (Al-Qadiri, Al-Alami, Al-Holy, & Rasco, 2008). DNA  
238oxidation and changes in DNA-protein interactions have been shown to be target sites for  
239oxidative reactions with NaOCl in the case of phage (Maillard, 1996). The PCA models for  
240NaOCl or PAA treated phage@anodisc are presented in **Figure 4**. The PCA results illustrated  
241that the spectral changes in the DNA region of a phage@anodisc is dose dependent and the PCA  
242model discrimination of spectral changes in NaOCl or PAA treated phage@anodisc was not  
243optimal. In the PCA model for PAA treated phage@anodisc, the PC1 and PC2 components  
244explained 78.9% and 13.1% of the variations in the spectral band corresponding to the DNA  
245region, respectively. In the PCA model of NaOCl treated phage@anodisc, the PC1 and PC2  
246components explained 77.4% and 16.5% of variations in the spectral band corresponding to the  
247DNA region, respectively. For PAA groups, the same PAA concentration corresponded to  
248mostly two clusters on the PCA visualization plot. For NaOCl, there was a high level of  
249overlapping between treatment dosages on the visualization plot. This was likely due to  
250nonlinear response of phage damage to NaOCl treatment. Various cellular responses to oxidation  
251reactions have been reported to have the nonlinear nature (Kalyanaraman et al., 2012; Neumaier  
252et al., 2012). In contrast, PCA is a linear transformation of the data. PCA analysis is generally  
253not recommended for modeling complex nonlinear relationship (Alanis-Lobato, Cannistraci,  
254Eriksson, Manica, & Ravasi, 2015).

255

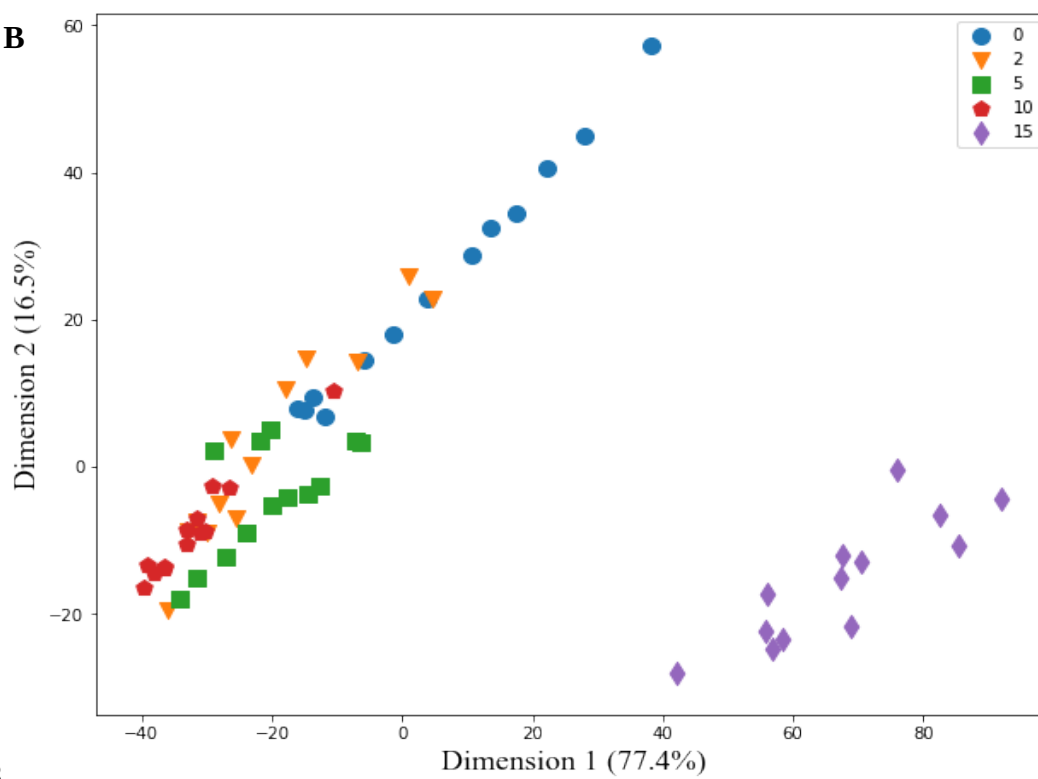
256

**A**



257

**B**



258

259

31

32



260 **Figure 4** 2D visualization of Principle Component Analysis of phage@anodisc FTIR  
261 spectra (A) PAA: 0, 20, 40, 60, 80 represents 0 ppm, 20 ppm, 40 ppm, 60 ppm, 80 ppm (B)  
262 NaOCl: 0, 2, 5, 10, 15 represents 0 ppm, 2 ppm, 5 ppm, 10 ppm, 15 ppm

263

## 264 Gradient Boosting model

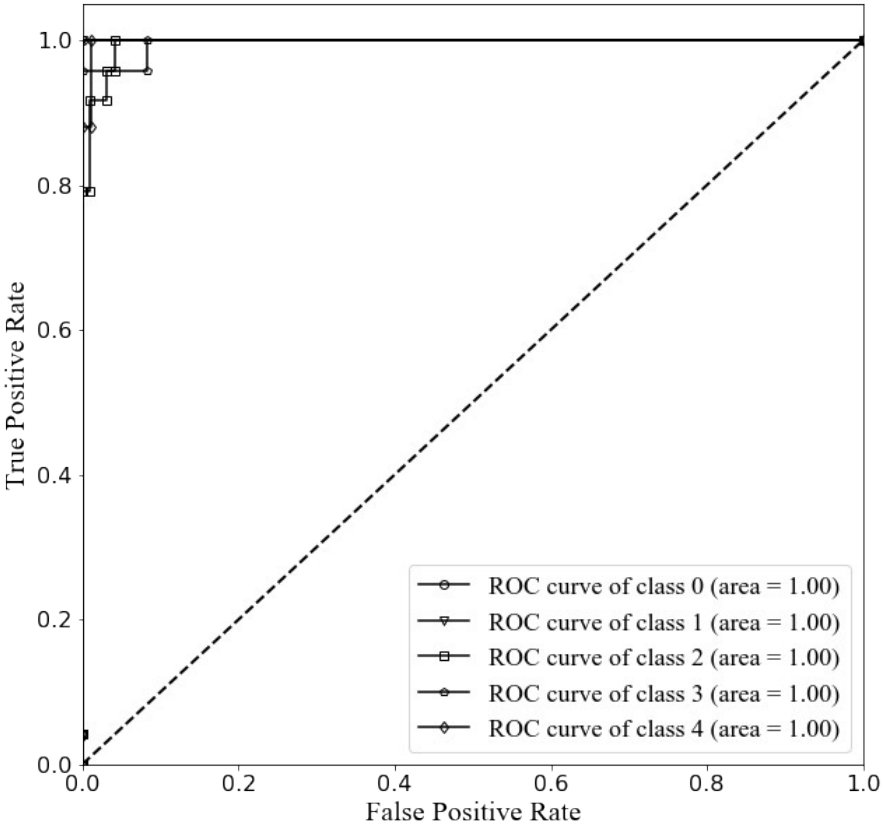
265 Due to nonlinear nature of the response observed based on the PCA analysis, the LightGBM was  
266 used for the analysis of the FTIR data. Using the LightGBM model, both sanitizer dosage as well  
267 as bacterial inactivation were predicted.

### 268 PAA Dosage

269 The prediction using the LightGBM model was evaluated based on the Receiver operating  
270 characteristic (ROC) curve and confusion matrix. A ROC curve is a plot of sensitivity on the y  
271 axis against (1-specificity) on the x axis for varying values of the threshold  $t$  (Zou, O'Malley, &  
272 Mauri, 2007). Sensitivity is defined as number of true positive samples (TP) / number of true  
273 positive or false negative (FP) samples. Specificity is defined as number of true negative samples  
274 / number of true negative or false positive samples. Threshold is the cut off probability for  
275 defining a positive class. The 45° diagonal line connecting (0,0) to (1,1) in the ROC curve  
276 corresponds to a random chance. The area under the ROC curve (AUC) is a summary measure  
277 that essentially averages diagnostic accuracy across the spectrum of test values. ROC is a  
278 suitable metric for balanced classification problem (Kotsiantis, Kanellopoulos, & Pintelas, 2006).  
279 The spectral dataset was balanced among all classes and hence ROC was selected as an  
280 evaluation matrix. Confusion matrix on the other hand provides a straightforward view of the  
281 number of data samples that have been correctly or incorrectly classified. **Figure 5** showed the

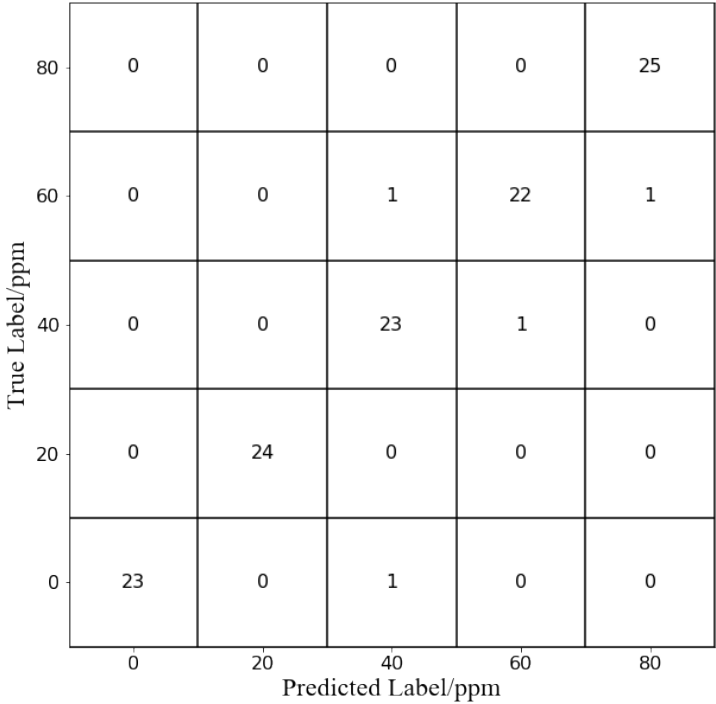
ROC and confusion matrix for predicting levels of PAA treatment, including 0 ppm (class 0), 20 ppm (class 1), 40 ppm (class 2), 60 ppm (class 3) and 80 ppm (class 4). In this result, for all 5 levels, ROC curves oriented towards the top left corner, indicating good prediction accuracy. Area under the curve (AUC) of ROC curve is another indicator for model performance. AUC for ROC curves among all levels reached 1, indicating effective classification of the spectral response. AUC results were consistent with the confusion Matrix. In the confusion matrix, classified samples are located in the diagonal part of the matrix. The total percentage of corrected predicted samples was 97% among a total of 121 samples. Data was also fitted with SVM model using Radial basis function kernel, but the prediction outcome was significantly worse (data not shown). In addition to the performance advantage of LightGBM, it is a very fast algorithm. The model training and prediction are completed within 3 mins for 121 data points. The LightGBM utilizes the Gradient-based One-Side Sampling and Exclusive Feature Bundling to expedite the calculation, which makes it suitable for handling big datasets common in industrial applications (Ke et al., 2017).

A



296

B



297

37

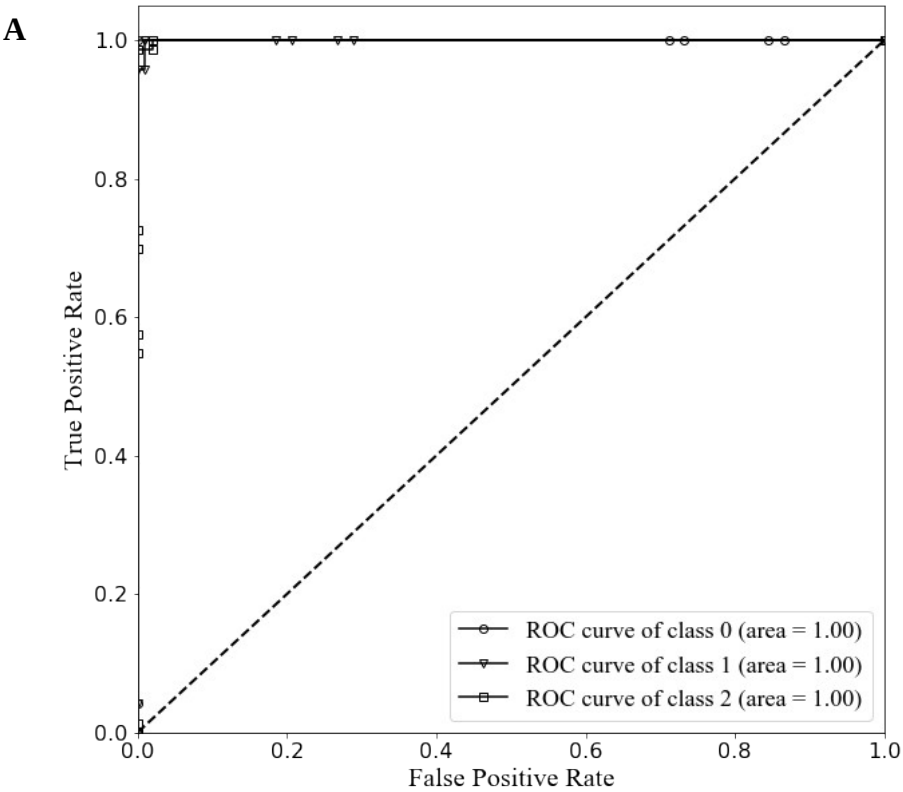
38

298**Figure 5** (A) Receiver operating characteristic plot (B) Confusion Matrix for prediction of PAA  
299concentration by phage@anodisc

### 300*Bacterial Inactivation by PAA*

301**Figure 6** showed the prediction of *E. coli* inactivation with FTIR data for PAA treatment. Data  
302was grouped into 4 classes based on the mean log reduction, namely 0 log reduction (class 0), 2  
303log reduction (class 1), >4 log reduction (class 2). The categorization was based on statistical  
304significant differences observed in the bacterial inactivation dataset. AUC for all classes were  
3051.0. According to the confusion matrix, the total percentage of correctly predicted samples was  
30696% among a total of 114 samples. Various studies have built machine learning models to  
307predict bacterial inactivation under sanitizers treatment including PAA using processing  
308conditions or bacterial internal signals as input variables. For example, Newhart *et al.* utilized  
309Artificial Neural Networks to model bacterial inactivation with physiochemical properties of  
310wastewater (Newhart et al., 2020). Caglar *et al.* predicted bacterial growth from mRNA and  
311protein abundances data (Caglar, Hockenberry, & Wilke, 2018). To the best of our knowledge,  
312this is the first study to predict bacterial inactivation with responses from surrogates using  
313spectroscopy methods.

314



**B**

class 2	0	0	73
class 1	0	21	3
class 0	23	0	1
	class 0	class 1	class 2

True Label

Predicted Label

319

320**Figure 6** (A) Receiver operating characteristic plot (B) Confusion matrix for prediction of *E. coli*  
321O157:H7 inactivation under PAA treatment

322

323*NaOCl Dosage*

324Hyperparameter tuning (HT) was utilized to improve prediction performance of the LightGBM  
325model for NaOCl as the default parameters did not obtain optimal outcome. Grid search  
326approach was used to conduct HT. Specifically, a set of hyperparameter combination was  
327selected as candidates. Models were trained based on all combinations in parallel. Model  
328performances were evaluated on test dataset. The hyperparameter used in the final model was  
329decided based on model performances. In the LightGBM, max tree depth, learning rate, sampling  
330rate and regularization terms are the major hyperparameters that affect the model performance  
331(Veronika, Vasily, & Kruchinin, 2018). **Table 1** showed the effect of Learning rate and  
332Regularization lambda L1 on the overall percentage of correct predicted samples. This metric  
333was chosen as it represents the overall model accuracy. Learning rate controls the step size at  
334each iteration when optimizing the objective function. Regularization lambda shrinks the model  
335coefficients terms to prevent overfitting. The choices of these two hyperparameters significantly  
336affected the total percentage of correct prediction. Based on this, learning rate 0.005 and Lambda  
33711 0.1 was chosen. Other hyperparameters and their values used in the study were as follows:  
338feature fraction (1), bagging frequency (1) and bagging fraction (0.7083). Feature fraction  
339defines the fraction of features to train each tree. Bagging frequency defines the frequency for  
340resampling from input dataset to build a tree whereas bagging fraction defines the fraction of  
341data to be used for each iteration. **Figure 7** showed the prediction of NaOCl concentration level  
342based on the phage FTIR data using optimal hyperparameters. ROC curve showed the model

43

22

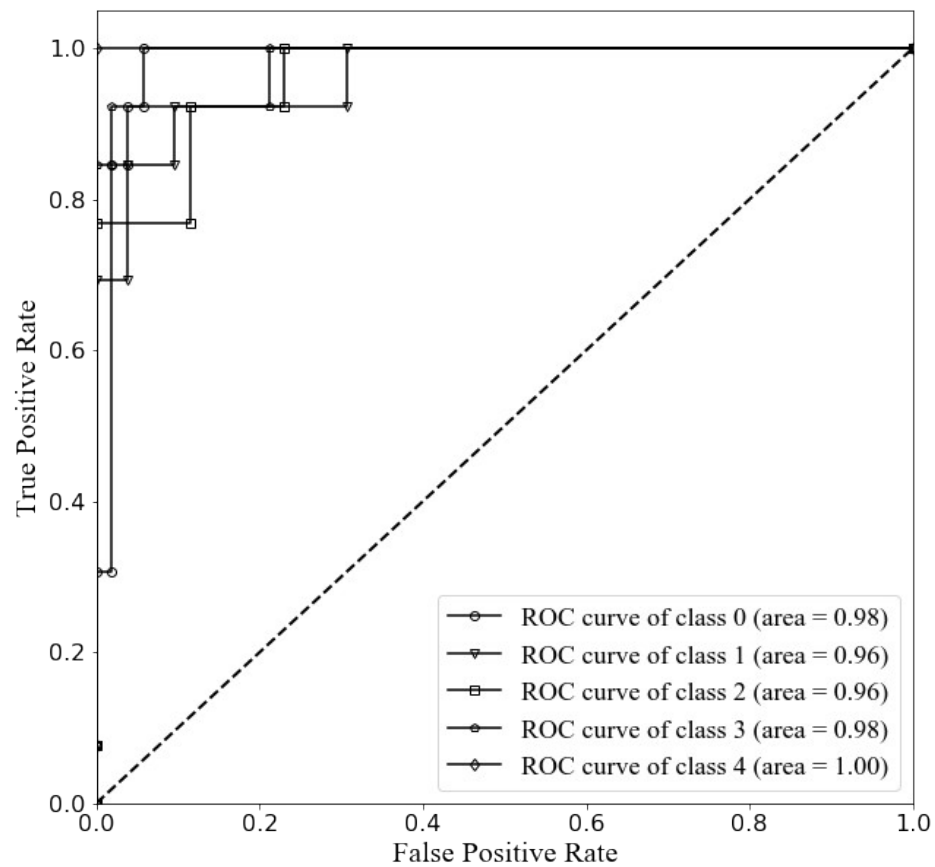
44

prediction performance has an AUC above 0.96 for all classes, namely 0 ppm (class 0), 2 ppm (class 1), 5 ppm (class 2) and 10 ppm (class 3) and 15 ppm (class 4). The confusion matrix also showed good prediction performance of the model. The total percentage of corrected classified samples was 88% among a total of 65 samples.

**Table 1** Hyperparameter tuning of LightGBM model for predicting NaOCl dosage

Learning rate	Lambda l1	Total corrected prediction percentage
0.0001	0.1	0.49
0.0001	0.3	0.32
0.0001	0.5	0.31
0.0005	0.1	0.82
0.0005	0.3	0.32
0.0005	0.5	0.31
0.005	0.1	0.88
0.005	0.3	0.32
0.005	0.5	0.31

**A**



**B**

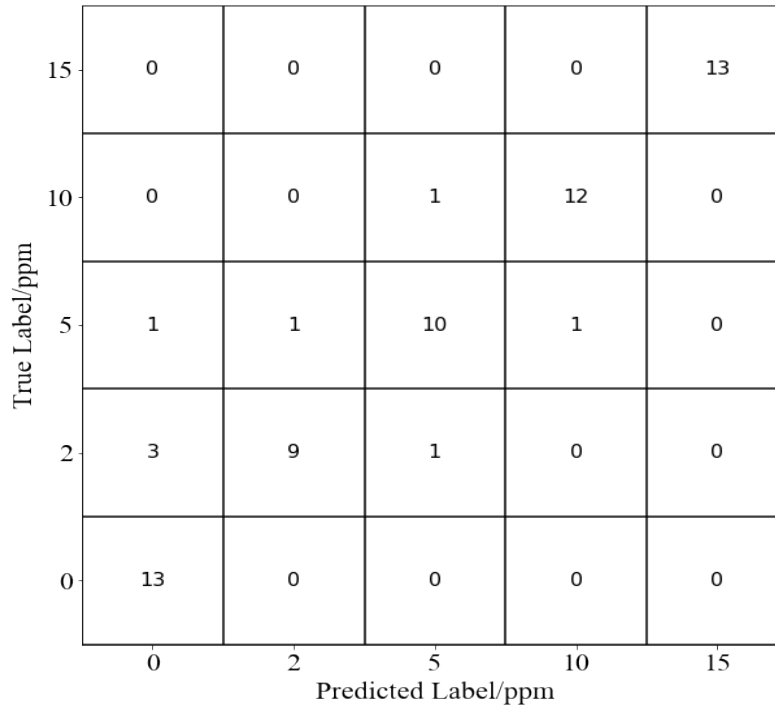
355

47

48

24





356

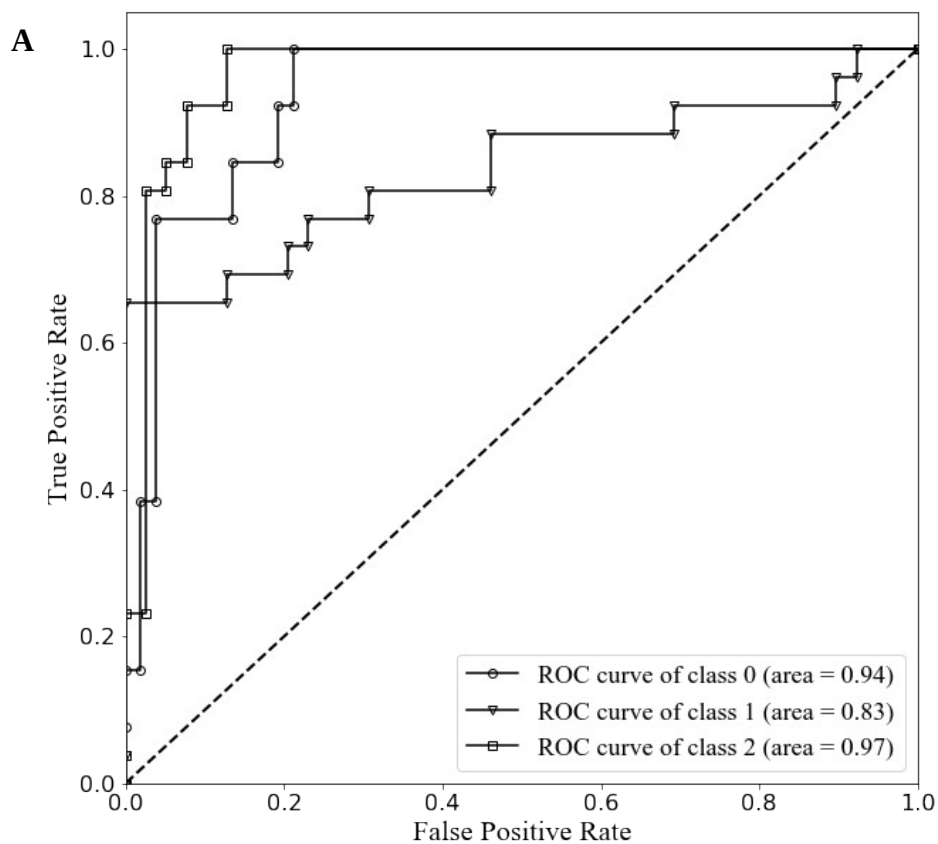
357**Figure 7** (A) Receiver operating characteristic plot (B) Confusion matrix for prediction of  
 358NaOCl concentration by phage@anodisc

359

360*Bacterial Inactivation by NaOCl*

361**Figure 8** showed the prediction of *E. coli* O157:H7 inactivation with FTIR data. *E. coli*  
 362inactivation was grouped into 3 classes, namely 0 log reduction (class 0), 3 log reduction (class  
 3631) and  $\geq 8$  log reduction (class 2). The categorization was also based on the statistical significant  
 364differences observed in the bacterial inactivation data with NaOCl treatment. AUC for all classes  
 365were at least 0.8. Among the three classes, class 2 had the smallest AUC (0.85), which could be  
 366attributed to large sampling variations within this class. According to the confusion matrix, the  
 367total percentage of true classification was 81% among a total of 65 samples. It is likely that the  
 368better model performance for predicting bacterial inactivation by PAA than NaOCl was due to a  
 369relatively large sample size in PAA. The performance of machine learning algorithm can be

highly impacted by the data size. Sordo *et al.* have shown that both Support Vector Machine (SVM) and Decision Trees show a substantial improvement in performance as the number of training samples increase (Sordo & Zeng, 2005).



**B**

True Label	class 2	3	4	19
	class 1	6	19	1
	class 0	11	1	1
		class 0	class 1	class 2
		Predicted Label		

375

376**Figure 8** (A) Receiver operating characteristic plot (B) Confusion matrix for prediction of *E. coli*  
377O157:H7 inactivation under NaOCl treatment

378

### 379Feature Importance Plot

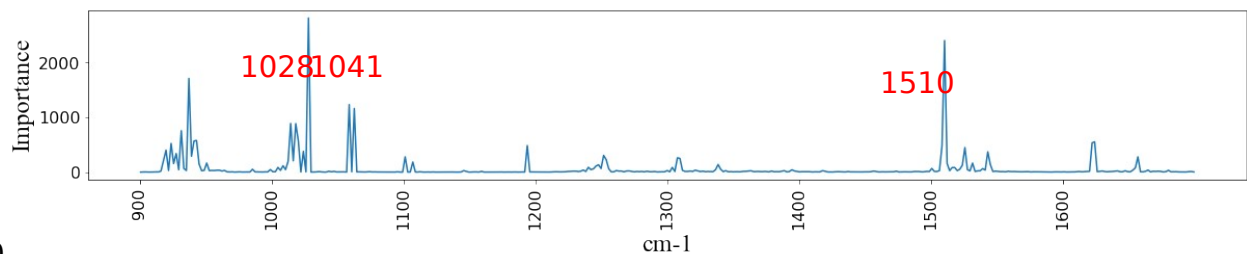
380Feature importance plot was used to reveal the contributions of each individual spectra  
381wavenumber to the prediction model. It was obtained by randomly shuffling each feature to  
382determine the increases in prediction error and then conducting ranking among all features  
383(Breiman, 2001). The feature importance of LightGBM model between 900 cm<sup>-1</sup> and 1700 cm<sup>-1</sup>  
384<sup>1</sup>was revealed as shown in **Figure 9**. The importance plot was extracted from the model based on  
385predicting bacterial inactivation under PAA treatment.

386

387

53

54



**Figure 9** Feature Importance plot for the prediction of *E. coli* O157:H7 inactivation under PAA treatment (x axis: wavenumber/cm<sup>-1</sup>; y axis: counts of the wavenumber being used to split the tree)

The x axis showed the spectra wavelength while the y axis showed the accumulated importance measure among all samples. Interestingly, the most important peaks were concentrated around 1000 cm<sup>-1</sup> as well as 1500 cm<sup>-1</sup>, which respectively correspond to the DNA and protein region in the FTIR spectra (Sahu et al., 2004; Simonova & Karamancheva, 2013). This coincided with the mechanisms of the biocidal action of PAA which involved both protein and DNA molecules (Alanis-Lobato et al., 2015; Kitis, 2004). Specifically, spectra at 1000 cm<sup>-1</sup> was assigned to

412 conformational changes in DNA due to single stranded DNA formation. The  $1083\text{ cm}^{-1}$  peak was  
413 assigned to symmetric phosphate groups in the DNA backbone also illustrating DNA  
414 fragmentation. Peaks at  $970$ ,  $1265$  and  $1041\text{ cm}^{-1}$  have been assigned to symmetric phosphate  
415 group, asymmetric phosphate group, stretching C-O ribose, and phosphate group, respectively  
416 indicating DNA fragmentation and changes in the deoxyribose structure (Oldenhof, Schütze,  
417 Wolkers, & Sieme, 2016; Ovissipour, Rai, & Nitin, 2019; Pascolo et al., 2016). The vibrational  
418 bands at  $1,044$  and  $1,113\text{ cm}^{-1}$  were also the characteristic markers of methionine oxidation  
419 (Ravi, Hills, Cerasoli, Rakowska, & Ryadnov, 2011). Amide I and amide II bands are the two  
420 major bands of the protein infrared spectrum. Peaks around protein amide I and II regions were  
421 also identified as key changes in the spectral features based on the prediction model. The protein  
422 amide I band (between  $1600$  and  $1700\text{ cm}^{-1}$ ) was related to the protein backbone conformation.  
423 Amide II (between  $1510$  and  $1580\text{ cm}^{-1}$ ) was associated with the N-H bending vibration and C-N  
424 stretching vibration (Barth, 2007). Peaks in this region showed high importance in predicting  
425 bacterial inactivation. Thus, the result indicated that both changes in structural conformation and  
426 oxidation related to DNA and protein were the key features used in the prediction model.

## 427 **Conclusions**

428 The overall goal of this study was to develop a rapid sanitation process verification based on  
429 spectroscopic measurement of immobilized phage oxidation and chemometric analysis to predict  
430 concentration of the selected sanitizers and bacterial inactivation. The results of this study  
431 indicate that vibrational spectroscopy coupled with machine learning models can be used for  
432 measuring, and quantifying phage responses to chlorine and PAA, and developing model for  
433 predicting the bacterial (*E. coli* O157:H7 as the reference) reduction, and sanitizers

434 concentrations. Overall, this study demonstrates immobilized phage as a surrogate for verifying  
435 the sanitation process in fresh produce industry.

436 **Acknowledgement:** This project was supported by grant no. 2015-68003- 23411 from the  
437 USDA-NIFA Program Enhancing Food Safety through Improved Processing Technologies  
438 (A4131) and by the USDA-AI Institute in Food Systems.

439

#### 440 **References**

441 Al-Qadiri, H. M., Al-Alami, N. I., Al-Holy, M. A., & Rasco, B. A. (2008). Using Fourier  
442 transform infrared (FT-IR) absorbance spectroscopy and multivariate analysis to study the  
443 effect of NaOCl-induced bacterial injury in water. *Journal of Agricultural and Food*  
444 *Chemistry*. <https://doi.org/10.1021/jf801604p>

445 Alanis-Lobato, G., Cannistraci, C. V., Eriksson, A., Manica, A., & Ravasi, T. (2015).  
446 Highlighting nonlinear patterns in population genetics datasets. *Scientific Reports*.  
447 <https://doi.org/10.1038/srep08140>

448 Barth, A. (2007). Infrared spectroscopy of proteins. *Biochimica et Biophysica Acta -*  
449 *Bioenergetics*. <https://doi.org/10.1016/j.bbabbio.2007.06.004>

450 Brackett, R. E., Ocasio, W., Waters, K., Barach, J., & Wan, J. (2014). Validation and  
451 verification: A practical, industry-driven framework developed to support the requirements  
452 of the food safety modernization Act (FSMA) of 2011. *Food Protection Trends*, 34, 410–  
453 425.

454 Breiman, L. (2001). Random forests. *Machine Learning*.

455 <https://doi.org/10.1023/A:1010933404324>

456 Caglar, M. U., Hockenberry, A. J., & Wilke, C. O. (2018). Predicting bacterial growth conditions  
457 from mRNA and protein abundances. *PLoS ONE*. <https://doi.org/10.1371/journal.pone.0206634>

458 Cossu, A., Le, P., Young, G. M., & Nitin, N. (2017). Assessment of sanitation efficacy against  
459 *Escherichia coli* O157: H7 by rapid measurement of intracellular oxidative stress,  
460 membrane damage or glucose active uptake. *Food Control*, 71, 293–300.

461 Gil, M. I., Selma, M. V., Lopez-Galvez, F., & Allende, A. (2009). Fresh-cut product sanitation  
462 and wash water disinfection: Problems and solutions. *International Journal of Food*  
463 *Microbiology*, 134, 37–45.

464 Kalyanaraman, B., Darley-USmar, V., Davies, K. J. A., Dennery, P. A., Forman, H. J., Grisham,  
465 M. B., ... Ischiropoulos, H. (2012). Measuring reactive oxygen and nitrogen species with  
466 fluorescent probes: Challenges and limitations. *Free Radical Biology and Medicine*.  
467 <https://doi.org/10.1016/j.freeradbiomed.2011.09.030>

468 Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., ... Liu, T. Y. (2017). LightGBM: A  
469 highly efficient gradient boosting decision tree. *Advances in Neural Information Processing*  
470 *Systems, 2017-December(Nips)*, 3147–3155.

471 Kitis, M. (2004). Disinfection of wastewater with peracetic acid: A review. *Environment*  
472 *International*. [https://doi.org/10.1016/S0160-4120\(03\)00147-8](https://doi.org/10.1016/S0160-4120(03)00147-8)

473 Kotsiantis, S., Kanellopoulos, D., & Pintelas, P. (2006). Handling imbalanced datasets : A  
474 review. *Science*.

475 Li, H., & Phung, D. (2014). Journal of Machine Learning Research: Preface. *Journal of Machine*

476 *Learning Research*, 39(2014), i–ii.

477 Maillard, J. Y. (1996). Damage to pseudomonas aeruginosa PAO1 bacteriophage F116 DNA by  
 478 biocides. *Journal of Applied Bacteriology*. [https://doi.org/10.1111/j.1365-](https://doi.org/10.1111/j.1365-2672.1996.tb03254.x)  
 479 [2672.1996.tb03254.x](https://doi.org/10.1111/j.1365-2672.1996.tb03254.x)

480 Morin, T., Martin, H., Soumet, C., Fresnel, R., Lamaudière, S., Le Sauvage, A. L., ... Maris, P.  
 481 (2015). Comparison of the virucidal efficacy of peracetic acid, potassium monopersulphate  
 482 and NaOCl on bacteriophages P001 and MS2. *Journal of Applied Microbiology*.  
 483 <https://doi.org/10.1111/jam.12870>

484

485 Neumaier, T., Swenson, J., Pham, C., Polyzos, A., Lo, A. T., Yang, P. A., ... Costes, S. V.  
 486 (2012). Evidence for formation of DNA repair centers and dose-response nonlinearity in  
 487 human cells. *Proceedings of the National Academy of Sciences of the United States of*  
 488 *America*. <https://doi.org/10.1073/pnas.1117849108>

489 Newhart, K. B., Goldman-Torres, J. E., Freedman, D. E., Wisdom, K. B., Hering, A. S., & Cath,  
 490 T. Y. (2020). Prediction of Peracetic Acid Disinfection Performance for Secondary  
 491 Municipal Wastewater Treatment Using Artificial Neural Networks. *ACS ES&T Water*.  
 492 <https://doi.org/10.1021/acsestwater.0c00095>

493 Oldenhof, H., Schütze, S., Wolkers, W. F., & Sieme, H. (2016). Fourier transform infrared  
 494 spectroscopic analysis of sperm chromatin structure and DNA stability. *Andrology*.  
 495 <https://doi.org/10.1111/andr.12166>

496 Ovissipour, M., Rai, R., & Nitin, N. (2019). DNA-based surrogate indicator for sanitation



497 verification and predict inactivation of *Escherichia coli* O157:H7 using vibrational  
 498 spectroscopy (FTIR). *Food Control*. <https://doi.org/10.1016/j.foodcont.2018.12.017>  
 499 Pascolo, L., Bedolla, D. E., Vaccari, L., Venturin, I., Cammisuli, F., Gianoncelli, A., ... Ricci, G.  
 500 (2016). Pitfalls and promises in FTIR spectromicroscopy analyses to monitor iron-mediated  
 501 DNA damage in sperm. *Reproductive Toxicology*.  
 502 <https://doi.org/10.1016/j.reprotox.2016.02.011>  
 503 Ravi, J., Hills, A. E., Cerasoli, E., Rakowska, P. D., & Ryadnov, M. G. (2011). FTIR markers of  
 504 methionine oxidation for early detection of oxidized protein therapeutics. *European*  
 505 *Biophysics Journal*. <https://doi.org/10.1007/s00249-010-0656-1>  
 506 Sahu, R. K., Argov, S., Salman, A., Huleihel, M., Grossman, N., Hammody, Z., ... Mordechai,  
 507 S. (2004). Characteristic absorbance of nucleic acids in the Mid-IR region as possible  
 508 common biomarkers for diagnosis of malignancy. *Technology in Cancer Research and*  
 509 *Treatment*. <https://doi.org/10.1177/153303460400300613>  
 510 Simonova, D., & Karamancheva, I. (2013). Application of Fourier transform infrared  
 511 spectroscopy for tumor diagnosis. *Biotechnology and Biotechnological Equipment*.  
 512 <https://doi.org/10.5504/BBEQ.2013.0106>  
 513 Sordo, M., & Zeng, Q. (2005). On sample size and classification accuracy: A performance  
 514 comparison. In *Lecture Notes in Computer Science (including subseries Lecture Notes in*  
 515 *Artificial Intelligence and Lecture Notes in Bioinformatics)*.  
 516 [https://doi.org/10.1007/11573067\\_20](https://doi.org/10.1007/11573067_20)  
 517 Suslow, T. (1997). *Postharvest Chlorination: Basic Properties & Key Points for Effective*

518     *Distribution. Postharvest Chlorination: Basic Properties & Key Points for Effective*  
519     *Distribution.* <https://doi.org/10.3733/ucanr.8003>

520 Veronika, A., Vasily, D., & Kruchinin, E. D. (2018). Why every GBDT speed benchmark is  
521 wrong. *ArXiv*.

522 Zareef, M., Chen, Q., Hassan, M. M., Arslan, M., Hashim, M. M., Ahmad, W., ... Agyekum, A.  
523     A. (2020). An Overview on the Applications of Typical Non-linear Algorithms Coupled  
524     With NIR Spectroscopy in Food Analysis. *Food Engineering Reviews*.  
525     <https://doi.org/10.1007/s12393-020-09210-7>

526 Zou, K. H., O'Malley, A. J., & Mauri, L. (2007). Receiver-operating characteristic analysis for  
527     evaluating diagnostic tests and predictive models. *Circulation*.  
528     <https://doi.org/10.1161/CIRCULATIONAHA.105.594929>

529

530

531

532

533

534

535