

Application of Convolutional Neural Networks for Prediction of 1 Disinfection By-Products

Nicolás M. Peleato (✉ nicolas.peleato@ubc.ca)

University of British Columbia

Research Article

Keywords: disinfection by-products, fluorescence spectroscopy, neural networks, machine learning, 22 water quality, natural organic matter

Posted Date: August 12th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-806261/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Scientific Reports on January 12th, 2022. See the published version at <https://doi.org/10.1038/s41598-021-03881-w>.

Application of Convolutional Neural Networks for Prediction of Disinfection By-Products

Nicolás M. Peleato^{1,*}

¹ University of British Columbia Okanagan, School of Engineering

* Corresponding author: nicolas.peleato@ubc.ca, 1137 Alumni Ave., Kelowna, British Columbia, V1V 1V7

Abstract

Fluorescence spectroscopy can provide high-level chemical characterization and quantification that is suitable for the use in on-line process monitoring and control. However, the high-dimensionality of excitation-emission matrices and superposition of underlying signals is a major challenge to implementation. Herein the use of Convolutional Neural Networks (CNNs) is investigated to interpret fluorescence spectra and predict the formation of disinfection by-products during drinking water treatment. Using deep CNNs, mean absolute prediction error on a test set of data for total trihalomethanes, total haloacetic acids, and the major individual species were all < 6 µg/L and represent a significant difference improved by 39% - 62% compared to dense neural networks. Heat maps that identify spectral areas of importance for prediction showed unique humic-like and protein-like regions for individual disinfection by-product species that can be used to validate models and provide insight into precursor characteristics. The use of fluorescence spectroscopy coupled with deep CNNs shows promise to be used for rapid estimation of DBP formation potentials without the need of extensive data pre-processing or dimensionality reduction. Knowledge of DBP formation potentials in near real-time can enable tighter treatment controls and management efforts to minimize the exposure of the public to DBPs.

Keywords: disinfection by-products; fluorescence spectroscopy; neural networks; machine learning; water quality; natural organic matter

1 Introduction

The use of fluorescence spectra for improved water quality monitoring¹ and as a process analytical technology for bioprocesses, food, and pharmaceutical production, has become increasingly popular.^{2,3} Fluorescence signatures are highly dependent on molecular structure, size, and environmental conditions, and therefore can be used to provide insight into chemical composition and properties.⁴ The sensitivity and specificity of fluorescence analysis, coupled with the potential real-time monitoring capabilities, fluorescence has applicability to a wide variety of process control applications.⁵

One promising application is the improved prediction and control of disinfection by-product (DBP) formation from drinking water treatment with chlorine. Chlorination is the most common disinfectant used worldwide. However, when chlorine reacts with natural organic matter (NOM), present in all natural water sources, various by-products of health concern are formed.⁶ Although many unique DBP species can be formed with varying public health risk, only specific groupings are commonly monitored and regulated in drinking water, including trihalomethanes (THMs) and haloacetic acids (HAAs). The monitoring frequency of regulated DBPs is generally low, with sampling only required once every 3 months for water systems in the United States and Canada.⁷

Significant attention has been paid to developing models that can predict DBP formation in order to improve knowledge of expected concentrations and inform water treatment operations of potential issues on a more frequent basis.⁸⁻¹² Since DBPs are formed from the reaction of chlorine and NOM, models must incorporate a measure of NOM. However, NOM is a chemically diverse grouping of organic

43 molecules whose characteristics are dependent on the surrounding environment. As such, the breadth of
44 potential NOM characteristics and the spatial and temporal variability results in significant challenges in
45 identifying an optimal measure that can capture this complexity and reactivity with chlorine.¹³
46 Fluorescence spectroscopy has considerable potential for the prediction and monitoring of DBP precursor
47 material. Many NOM compounds fluoresce and fluorescence measures can capture some chemical
48 characteristics of NOM.¹⁴ Previously, fluorescence has been used with success to predict or identify
49 correlations with regulated DBPs,¹⁵⁻¹⁷ as well as unregulated or by-products of emerging concern such as
50 chloral hydrate¹⁸ and haloacetonitriles.^{19,20}

51 A common challenge to implementing fluorescence as a monitoring tool is the high-dimensionality
52 and superposition of the resulting emissions. When utilizing fluorescence spectra collected at iterated
53 excitation/emission wavelengths, a dimensionality reduction approach is often used to simplify excitation-
54 emission matrices (EEMs).¹⁴ By identifying a few underlying components that explain most of the variance
55 in the data, the hypothesis is that noise is reduced, and subsequent modelling using a reduced
56 dimensionality improves prediction. A basic simplification or dimensionality reduction approach would be
57 to select peaks or regions in the fluorescence spectra where regional integration or peak fluorescence can
58 be determined. While this type of expert guided approach has been used extensively in the past,
59 discarding the majority of collected data neglects the richness of information contained. For complex
60 systems such as those that include identifying natural organic matter (NOM) in water, organic
61 fluorophores with similar chemical structures are not easily distinguished in the spectra. The use of
62 principal component analysis (PCA) or parallel factors analysis (PARAFAC) has revealed underlying signals
63 resembling fluorophores, which can be tied to spectral regions from which chemical properties can be
64 inferred.^{14,21,22} These analysis approaches are often limited to linear dimensionality reduction, so non-
65 linear features such as Rayleigh or Raman scattering need to be removed from the spectra.²³
66 Furthermore, potential impacts of environmental conditions such as pH or temperature,⁴ or possible
67 charge-transfer interactions²⁴ may invalidate the assumption of a linear relationship between fluorophore
68 concentrations and fluorescence intensity. Inner filter effects are also prevalent, where incident excitation
69 light and emitted fluorescence is quenched by other chromophores present in the sample, result in a non-
70 linear intensity response.²⁵ Constraints imposed by the method can be helpful when derived from prior
71 knowledge of the system, such as non-negativity of fluorescence emissions, reducing bias and possibly
72 resulting in a more accurate depiction of underlying structures. However, these same constraints may
73 limit the overall accuracy of reconstruction based on the condensed representation.²⁶

74 It may be advantageous to directly use all data collected in fluorescence EEMs to limit potential errors
75 introduced from dimensionality reduction. However, for water quality analysis, there have been limited
76 studies that explore the use of full fluorescence EEMs without dimensionality reduction. Non-linear
77 regression using high-dimensional inputs can be accomplished using neural networks. More recent work
78 with Convolutional Neural Networks (CNNs or ConvNets) has shown this type of network structure is well
79 suited to interpreting images or other tasks datasets with local groups of values that are highly
80 correlated.²⁷ Instead of training weighted connections between every individual node, CNNs train spatial
81 filters or kernels to identify small recurring features in the input space. The use of filters allows for
82 parameter sharing, where trained weights are used throughout the input space and are not tied to specific
83 input nodes, giving rise to spatial invariance of features.²⁸ Furthermore, CNNs typically employ pooling
84 layers where outputs in specific locations are merged with nearby outputs, creating invariance to small
85 distortions in the input and reducing the dimensionality of the representation.^{27,28} CNNs have been
86 successfully applied in chemometric applications such as interpreting Raman and mid-infrared spectra for
87 identifying *Escherichia coli* and meats,²⁹ pharmaceuticals in tablets with near infrared spectra,³⁰
88 categorizing wines using infrared spectra,³¹ and classification of manganese valence.³² However, there has
89 been no use of CNNs for interpreting 2D fluorescence spectra, and previous implementations have
90 focused on 1D infrared or Raman spectra. Furthermore, the use of CNNs for fluorescence analysis of water

91 quality has not been explored. It is hypothesized that the strengths of CNNs for processing and
92 interpreting spatially dependent data will be well suited for 2D fluorescence spectra where local groups
93 of values are highly correlated.

94 This paper investigates the use of deep NNs and CNNs to interpret fluorescence spectra for the
95 prediction of DBP formation potential. The two major groups of regulated DBPs are assessed, THMs and
96 HAAs including the individual species that made up these groups in the samples analysed
97 (trichloromethane, bromodichloromethane, trichloroacetic acid, and dichloroacetic acid). A method to
98 interpret the CNN results is also used to identify fluorescence regions that are most likely associated with
99 high DBP formation potentials.

100 2 Results

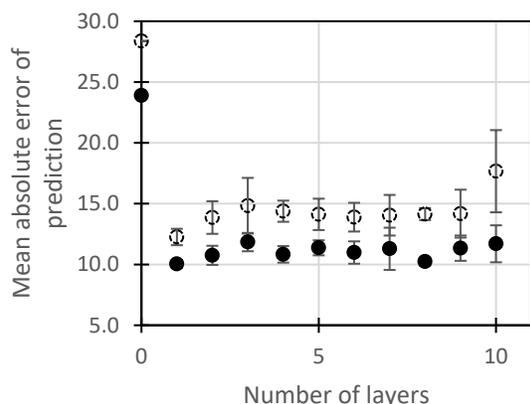
101 A dataset of DBP formation potentials and associated fluorescence EEM measurements were used to
102 assess the capabilities of deep NNs and CNNs for water quality analysis. Water samples analyzed were
103 from a pilot-scale treatment plant receiving river water. Samples were taken throughout a treatment train
104 consisting of several unit processes including coagulation, flocculation, sedimentation, ozonation,
105 advanced oxidation (peroxide and ozone), and filtration through anthracite or activated carbon. As such,
106 the samples analyzed had a wide range of NOM concentrations and characteristics. Dissolved organic
107 carbon varied from 2.6 to 6.3 mg L⁻¹, and specific ultraviolet absorbance varied from 0.75 to 2.53 L mg⁻¹
108 m⁻¹ over all samples. DBP formation potentials were determined by maintaining a free chlorine residual
109 of 1.5 mg/L for 24 hours. Although all four chlorinated or brominated THM and nine HAA species could be
110 detected, only trichloromethane (TCM), bromodichloromethane (BDCM), trichloroacetic acid (TCAA), and
111 dichloroacetic acid (DCAA) were consistently identified at concentrations above detection limits.

112 2.1 Dense Networks

113 An iterative optimization approach was used to understand the impact of NN structure on overall
114 performance. While many aspects of network structure can be optimized, the focus was on the number
115 of layers (i.e. depth). A dense network was trained with an increasing number of layers to identify the
116 degree to which network depth can improve prediction accuracy on a test set. Figure 1 shows the total
117 THM and HAA predictions results given the number of layers in a dense NN. The error bars in Figure 1
118 represent the standard deviation of 8 repeated random initializations of the network. A network with 0
119 hidden layers is simply the input values (dimensions = 6,336) connected to 1 output node. When the
120 number of layers was increased, each layer's nodes were set to half of the previous layer. For example,
121 with two hidden layers, hidden layer 1 would have 3,168 nodes, and layer 2 would have 1,584 nodes.

122 As observed in Figure 1, increasing the number of hidden layers beyond one in a dense NN did not
123 improve network performance in predicting both total THMs and HAAs. Total THM mean absolute error
124 was at a minimum with one hidden layer (MAE: 12.26 ± 0.34 µg/L) and total HAA error was also at a
125 minimum with one hidden layer (MAE: 10.04 ± 0.33 µg/L). Similar results were observed for individual
126 species (Table 2). While adding additional layers eventually reduced test set error, this came at the cost
127 of increased variability between random network initializations. For example, the coefficient of variation
128 (CV) increased from 3.3% for a one-layer dense network to predict total HAAs to 8.4% for 6 layers.
129 Increased variability in performance could be due to the increased number of learnable parameters and
130 the relatively small sample size used in this study. The structure chosen, where each layer contained half
131 the nodes of the previous layer, resulted in 15,870,977 trainable parameters with one hidden layer.
132 However, it should be noted that a dense network with one hidden layer of 10 nodes (56,361 trainable
133 parameters) resulted in a decrease in prediction accuracy (1 layer, 10 nodes: total THM MAE 13.54 ± 0.47
134 µg/L; total HAA MAE 12.62 ± 3.00 µg/L). The lack of improvement of dense NNs beyond one hidden layer

135 is expected given the small data size and demonstrates that deep dense networks are unlikely to provide
136 advantages in modelling small water quality datasets.



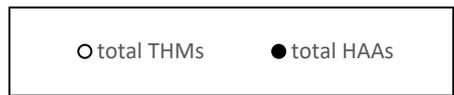
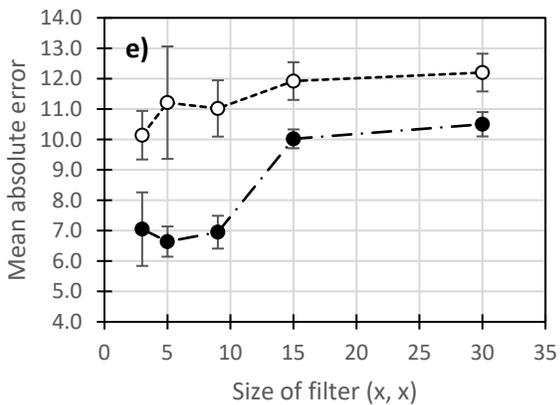
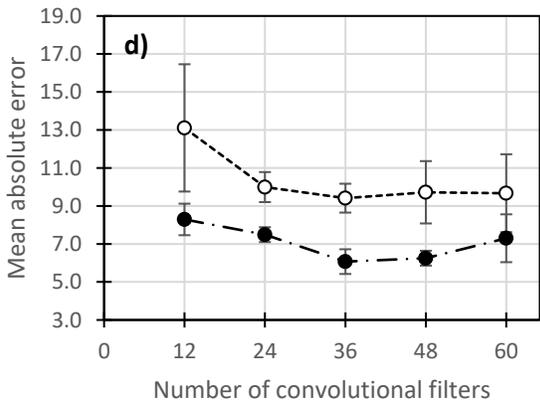
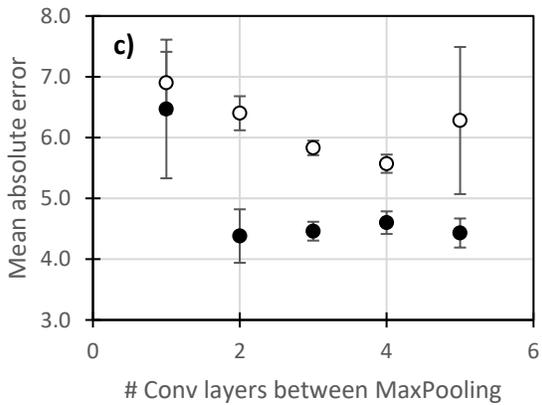
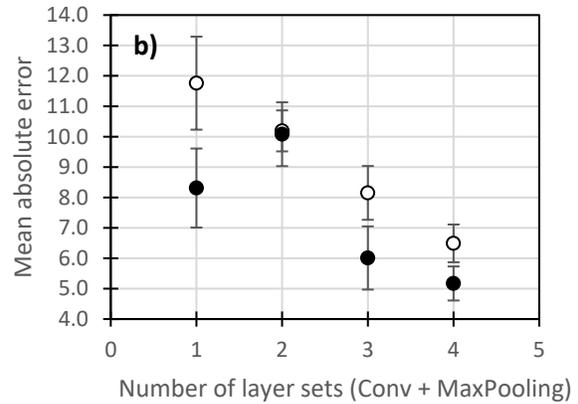
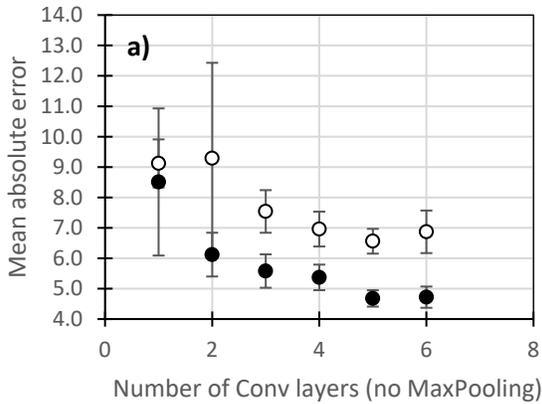
137 **Figure 1** Mean absolute error (MAE) of prediction for a test set ($n = 28$) of a) THM concentrations and b)
138 HAA concentrations using a dense neural network structure. Error bars represent the standard deviation
139 of MAE over 8 random initializations of the network weights.
140

141 2.2 Convolutional Networks

142 In contrast to dense NNs, prediction accuracy was minimized with increasing CNN network depth (Figure
143 2). Network depth was investigated by increasing the number of convolutional layers and the number of
144 layer sets (convolution followed by max pooling). Convolutional layers provide learned filters or kernels
145 that identify small features in the spectrum, while max pooling layers decrease dimensionality and pool
146 redundant features.^{27,28} Including 5 hidden convolutional layers without pooling layers was found to
147 optimize prediction accuracy (Figure 2a). Increasing the number of layer sets also significantly improved
148 performance (Figure 2b). A further decrease in error of 19.3% for total THMs and 28.9% for total HAAs
149 was observed when increasing both the number of convolutional layers between max pooling layers and
150 increasing the number of max pooling layers (Figure 2c; Table 2).

151 It was also of interest to investigate the role of the size of receptive fields for each filter and the
152 number of filters included in each convolutional layer. The receptive field size identifies the number of
153 adjacent data points to be considered by each filter. Previous work in chemometrics has shown relatively
154 large receptive fields to work well³¹ and could identify features that span over large areas of the spectra,
155 however, expanding the filter size increases the number of trainable parameters. Alternatively, by
156 including convolutional layers in sequence, the receptive field's effective size is expanded, minimizing the
157 number of trainable parameters and including additional layers of non-linearity.³³ As such, the results
158 suggest that larger receptive fields could improve CNN performance. However, increasing the receptive
159 field size beyond (3, 3) for individual layers did not improve performance (Figure 2e), and expanding
160 receptive fields may be best accomplished by stacking convolutional layers in sequence. It is also of note
161 that increasing the number of trained filters improved performance up to 36, after which no further
162 changes were observed (Figure 2d).

163 An example of the learned filters are shown in Figure 3. CNNs create hierarchical representations of
164 data showing how specific irrelevant spectral features are discarded and specific areas of the spectra
165 needed to predict DBP concentrations are magnified.²⁹ The initial filter layer identifies large and smooth
166 and broad features in the spectra. After pooling, feature maps become more coarse and more distinct
167 patterns between filters can be discerned, highlighting specific areas of the spectra. In the last layer of
168 feature maps, many filters highlight one constant emission level over several excitation bands (left to
169 right).



170

171

172

173

174

175

176

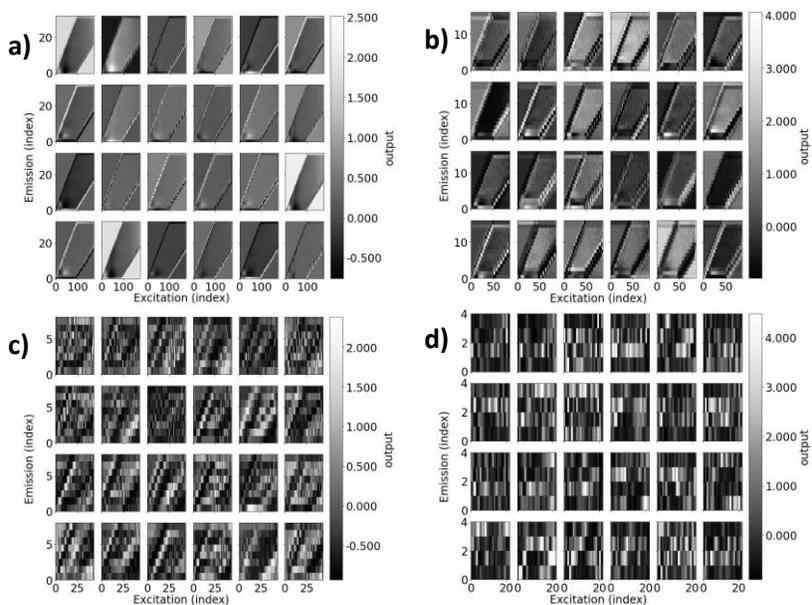
177

178

179

Figure 2 Impact of CNN structure and depth on MAE of test set predictions. a) varies the number of convolutional layers without any max pooling layers, b) varies the number of layer sets with 1 convolutional layer followed by max pooling, c) varies the number of convolutional layers between max pooling layers (4 max pooling layers in total), d) varies the number of convolutional filters for 1 convolutional layer without max pooling, e) varies the size of the receptive field for 1 convolutional layer without max pooling.

180



181

182

183 **Figure 3** Feature maps of convolutional filters (24) from 4 convolutional layers chosen between max
 184 pooling layers. a) first convolutional layer, b) after the first max pooling, c) after the second max pooling,
 185 d) after the third max pooling. All max pooling was carried out over a (2,2) window.

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

Table 1 Mean absolute error (MAE) of predictions on tests set for several model types. Range or error is the standard deviation from 8 repeated initializations of the same model (where applicable). Error (\pm) is calculated as the standard deviation over 8 random initialization of the network weights.

Disinfection By-Product Species	Model			
	Dense (1 layer)	CNN (1 layer)	CNN (4 pooling layers, 1 convolutional layer)	CNN (4 pooling layers, 4-5 convolutional layers)
Total THMs	12.26 \pm 0.34	6.62 \pm 0.13	6.06 \pm 0.22	5.57 \pm 0.15
Trichloromethane (TCM)	8.80 \pm 0.32	7.00 \pm 0.84	4.87 \pm 0.62	3.39 \pm 0.10
Bromodichloromethane (BDCM)	6.34 \pm 0.51	4.69 \pm 0.33	4.11 \pm 0.22	3.86 \pm 0.06
Total HAAs	10.04 \pm 0.33	4.54 \pm 0.24	4.20 \pm 0.47	4.43 \pm 0.13
Dichloroacetic acid (DCAA)	7.82 \pm 0.43	6.07 \pm 1.01	4.84 \pm 0.26	4.19 \pm 0.08
Trichloroacetic acid (TCAA)	8.36 \pm 0.33	5.85 \pm 0.90	4.58 \pm 0.16	4.22 \pm 0.10

202 2.3 Model explanations

203 The objective of identifying model explanations was to confirm that that model predicts high or low
204 concentrations of DBPs based on fluorescence features that are known or possibly associated with DBP
205 precursors. There are scattering signals (i.e. not from organic material) or other potential artifacts from
206 the sample analysis process that would bias the model to “know” concentrations of DBPs for incorrect
207 reasons. The second objective was to identify fluorescence regions most highly associated with specific
208 DBP formation potentials. This information could be used to further understanding of the characteristics
209 of DBP precursors and potentially optimize treatment processes that preferentially remove compounds
210 with those characteristics.

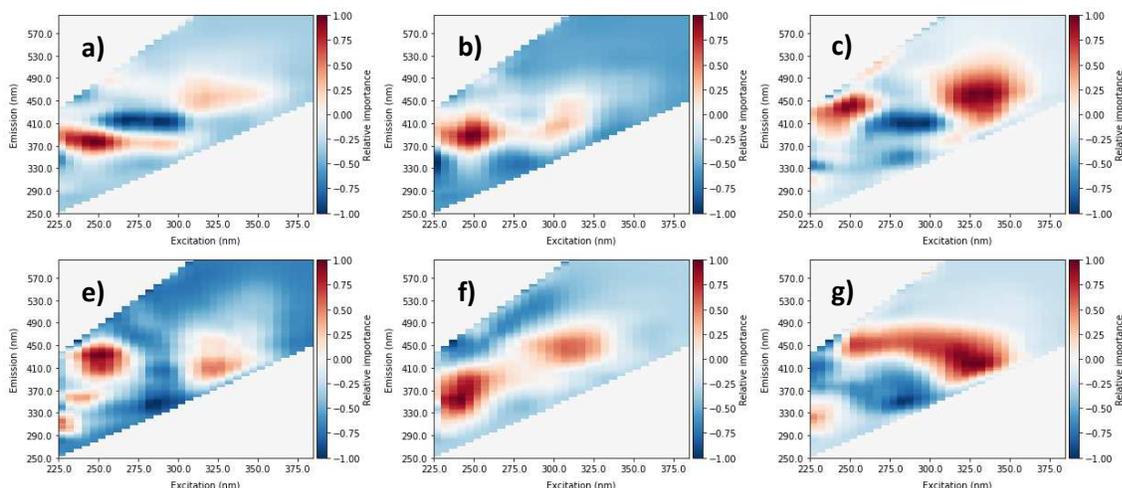
211 An occlusion method was applied to identify spectral areas that most significantly influence prediction
212 accuracy. The occlusion method identifies spectral regions most relevant to a prediction by randomly
213 occluding or setting a segment of all inputs in a specified region to 0. The error incurred due to this
214 occlusion indicates how relatively important that specific area is to accurate predictions. The error was
215 calculated as the difference between non-occluded and occluded predictions, and the direction or sign of
216 the error was preserved. As such, positive values indicate that the model underestimated DBP formation
217 with a specific patch occluded, and negative values indicate overestimated DBP formation. A total of
218 20,000 iterations of random patches per model were chosen to build the heat maps. A random approach
219 to selecting the patch was taken to reduce any bias from neighbouring values since the variables included
220 in each patch would change between iterations. Figure 4 shows the average heat maps identified from
221 training deep CNNs on total DBPs and individual species. Likewise, Figure 5 shows heat maps based on
222 dense NNs.

223 From occlusion heat maps of variable or spectra area importance (Figure 4), it was observed that
224 fluorescence in the area of approximately ex: 225 – 260 nm and em: 370 – 500 nm was most impactful of
225 prediction accuracy for all DBPs (both total and individual species). A second common area of importance
226 at ex > 300 nm and em > 400 nm was also observed. Fluorescence in these two regions is generally
227 considered to be humic-like and fulvic-like material.³⁴ Several heat maps also show areas of importance
228 in protein-like fluorescence regions (excitation: 230 – 250 nm, emission: 300 – 360 nm) associated with
229 tryptophan-like or tyrosine-like fluorescence.³⁴ Spectral importance in these regions conforms well to
230 expectations of DBP precursor type material that can fluoresce, generally thought to be aromatic humic-
231 like or fulvic-like material.³⁴ The heatmaps provide evidence that the NNs are utilizing signals from regions
232 that are reasonable for DBP prediction. Previous DBP prediction methods based on fluorescence data have
233 utilized the same spectral regions.^{15,16,35}

234 Compared to dense NNs, CNN heat maps show broader areas of importance with more gradual
235 changes (Figure 4 and Figure 5). Gradual changes conform with the expectation of fluorescence signals
236 from fluorophores, and sharp changes are not typically associated with fluorescence from natural organic
237 matter.¹⁴ Furthermore, CNN heat maps emphasize higher excitation bands. Particularly for prediction of
238 trihalomethanes, peaks at excitation > 300 nm were important for positive predictions, while dense NN
239 heat maps placed less importance on these areas. Several CNN heat maps show signals that have several
240 excitation peaks, with limited changes in emission. For example, the BDCM heat map shows peaks at
241 approximately 250 nm and 330 nm, with emissions constant at 450 nm. Multiple excitation peaks at one
242 singular emission conform well with expectations of fluorescence from individual fluorophores, where
243 multiple wavelengths can cause excitation, however, the emission is always from the lowest singlet state
244 and therefore only at one wavelength.³⁶ The identification of areas of importance at several emission
245 bands suggests several distinct fluorophores contributing to DBP formation potential rather than
246 individual components.

247
248

249



250

251

252

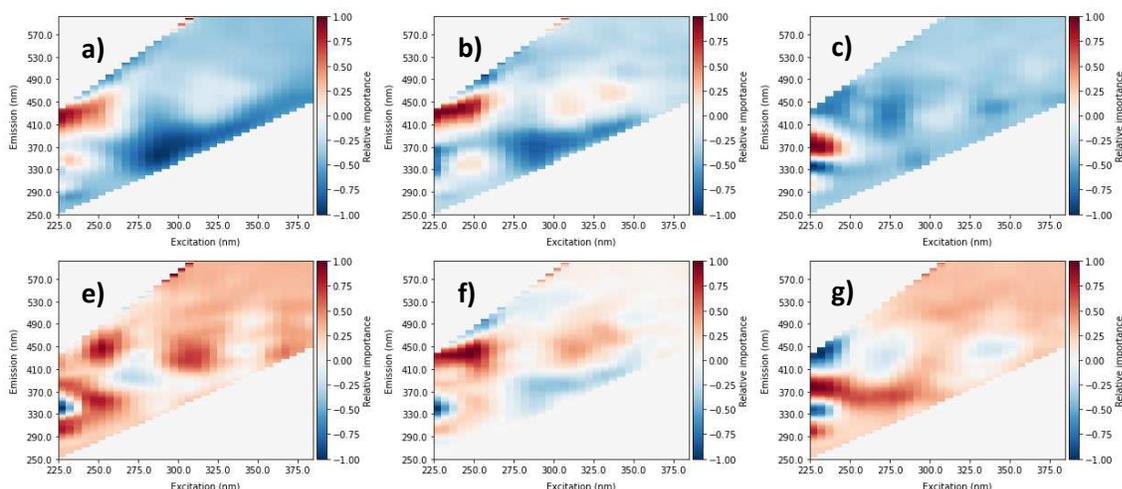
253

254

255

Figure 4 Heat maps from random occlusion of variable importance for CNN prediction of a) total THMs, b) trichloromethane, c) bromodichloromethane, d) total HAAs, e) trichloroacetic acid, f) dichloroacetic acid.

256



257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

Figure 5 Heat maps from random occlusion of variable importance for dense network prediction of a) total THMs, b) trichloromethane, c) bromodichloromethane, d) total HAAs, e) trichloroacetic acid, f) dichloroacetic acid.

It can also be observed that there is greater continuity of areas of importance with individual DBP species and the total levels. Since the individual species should sum to total concentrations, total THM and total HAA heat maps would be expected to show similar characteristics to the individual species. Ex/em 350 nm/380 nm is observed in CNN heat maps for total THMs (Figure 3a) and TCM (Figure 3b). The secondary peak for total THMs at approximately ex/em 325 nm/450 nm is mirrored in the BDCM heat map (Figure 3c). Similar conformance was not observed with dense NN heat maps, for example, the areas of highest importance for BDCM (Figure 4c) was not present in the total THM heat map (Figure 4a). However, while some overlap is present between CNN heat maps of species and the total DBP levels, it should be noted that not all peaks are mirrored (e.g. BDCM peak at ex/em 250 nm/450 nm not seen in total THM map).

272 Individual species of DBPs showed differences between CNN heat maps. BDCM areas of importance
273 were shifted to higher excitation and emission areas compared to TCM. A similar pattern can be seen
274 between DCAA and TCAA. Identified differences in spectral areas between individual species were
275 expected since preferential yields of specific by-products from pure model compounds have pointed to
276 certain molecular structures resulting in the preferential formation of individual DBP species.^{37,38} A shift
277 to the greater importance of fluorophores at emissions > 450 nm could indicate BDCM and DCAA
278 formation resulting from humic-like material with greater oxygen/carbon ratios and lower
279 hydrogen/carbon ratios, implying an oxidation state ≥ 0 .³⁹

280 A second notable difference is the increased importance of protein-like material (ex/em 230 – 250
281 nm/300 – 350 nm) for HAA predictions. This peak location is typically associated with aromatic amino
282 acids such as tryptophan and tyrosine.³⁴ Previous studies show that aromatic amino acids⁴⁰ and protein-
283 like fluorescence signals strongly correlate with HAA formation potentials.^{19,41,42} In particular, the protein-
284 like peak was observed to be most prominent for the prediction of TCAA. This observation conforms well
285 to previous results that show higher TCAA formation than DCAA from aromatic amino acids that would
286 contribute to the observed fluorescence signal.⁴⁰ From the dense NN heat map of TCAA, regions
287 surrounding the expected aromatic amino acid peak are positive. However, there is a strong negative
288 relationship in the specific location of tryptophan fluorescence (ex/em 230 nm/340 nm).

289 3 Discussion

290 This study investigated the use of deep CNNs to interpret fluorescence spectra and predict the formation
291 of regulated chlorination DBPs from a drinking water treatment plant. The observed results indicate that
292 deep CNNs are well suited to the task of interpreting fluorescence excitation-emission matrices and
293 prediction of DBPs for several reasons: 1) overall prediction accuracy for all DBP groups and species were
294 significantly reduced compared to dense NNs and previous modelling approaches using dimensionality
295 reduction, 2) results from random initializations were less variable using deep CNNs compared to dense
296 and shallow CNNs, 3) deep CNN heat maps show trained networks utilize data from spectral regions that
297 are well known to be associated with DBP formation potentials, and 4) compared to dense NNs, CNNs
298 show heat maps with characteristics more conformant with expectations of fluorescence from organic
299 precursor material.

300 Compared to previous work that utilized dimensionality reduction prior to regression, the use of CNNs
301 significantly improved the accuracy of prediction. Using the same dataset and training/test data, optimal
302 results of total THM MAE 7.46 $\mu\text{g/L}$ and total HAA MAE 10.75 $\mu\text{g/L}$ were previously reported.¹⁷ Significant
303 reduction of prediction accuracy was achieved using CNN architectures, particularly for HAA prediction
304 (this study: total THM MAE 5.53 $\mu\text{g/L}$ and total HAA MAE 4.37 $\mu\text{g/L}$). It is not straightforward to compare
305 results to other studies given the variation in number of samples, methods for formation potential
306 determination, range of concentrations in the training/test sets, and performance metrics. However,
307 results found in this study also represent improvements over HAA and THM formation predictions
308 previously reported using similar performance metrics (e.g. total THM MAE 13.5 $\mu\text{g/L}$ and total HAA MAE
309 7.7 $\mu\text{g/L}$).¹⁵ Furthermore, previous approaches to utilize fluorescence data for DBP predictions have relied
310 heavily on dimensionality reduction to identify relevant fluorescence features, which adds complexity and
311 may neglect to capture features of importance. In contrast, deep CNNs present an opportunity to utilize
312 full fluorescence spectra without the need for manual or highly supervised feature selection through peak-
313 picking, regional integrations, or PARAFAC analysis. It is thought that deep CNNs provide an opportunity
314 for complex behaviours to be represented by several simpler representations. Observation of feature
315 maps produced by convolutional layers show a hierarchy of feature representations, with general and
316 smooth representations in high layers and progressively coarser and more specific highlighted spectral
317 areas as increasing numbers of convolution and pooling layers are applied.

318 Neural networks are often discussed as black-box type algorithms, where the internal reasoning is
319 unknown or difficult to illustrate. However, it is imperative that the logic of prediction algorithms used in
320 applied tasks, such as prediction of potentially toxic disinfection by-products, is discernible. There is also
321 an opportunity to use these powerful data-driven approaches to help identify important variables or
322 characteristics of the system. The use of heat maps generated from an occlusion approach to identifying
323 spectral areas that highly influence predictions gives insight into the decision-making process and helps
324 confirm that trained networks are relying on data from spectral regions associated with DBP precursors.
325 Furthermore, heat maps can help direct future more detailed studies investigating the characteristics of
326 precursor material.

327 As such, the use of fluorescence spectroscopy coupled with machine learning techniques, such as deep
328 CNNs, show promise to be used for rapid estimation of DBP formation potentials. In the context of typical
329 regulatory thresholds for water treatment (total THMs < 80 µg/L, total HAAs < 60 µg/L) the presented
330 methodology produced error levels (MAE 3.39 – 5.53 µg/L) that would be appropriate for rapidly
331 informing operations and management regarding conformance with regulatory thresholds. Knowledge of
332 DBP formation potentials in near real-time can enable tighter treatment controls and management efforts
333 to minimize the public's exposure to DBPs.

334 4 Methods

335 4.1 Water samples

336 Water samples were obtained from parallel pilot treatment trains that were fed Otonabee River water
337 (Peterborough, Ontario, Canada). Samples were obtained throughout the treatment train for fluorescence
338 analysis and for determining DBP formation potentials. Processes applied included coagulation,
339 flocculation, sedimentation, ozonation, advanced oxidation (peroxide and ozone), and filtration through
340 anthracite or activated carbon. Further information on the pilot-scale set-up and water samples can be
341 found in Peleato et al. (2017).⁴³

342 4.2 Fluorescence

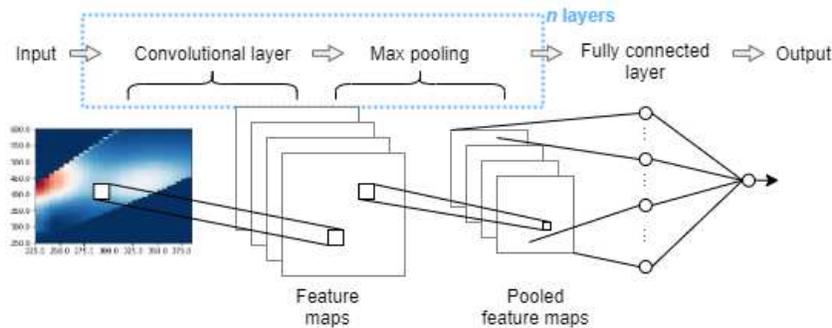
343 Fluorescence spectra were collected using an Agilent Cary Eclipse fluorescence spectrophotometer
344 (Mississauga, Canada). Excitation and emission wavelength ranges were 225 – 380 nm (5 nm increments),
345 and 250 – 600 nm (2 nm increments), respectively. The fluorescence spectra were blank subtracted using
346 Milli-Q® water. The spectrum for Milli-Q® water was also used to apply Raman corrections at an excitation
347 wavelength of 350 nm and bandwidth of 5 nm to allow fluorescence intensities to be reported in Raman
348 Units (RU).⁴⁴ Absorbance values collected over the excitation-emission range were used to correct for
349 inner-filter effects. Rayleigh scattering lines were removed by setting all values above 2nd order Rayleigh
350 or below 1st order Rayleigh to 0. The absorbance corrected spectra were then scaled between 0 and 1 for
351 each excitation/emission pair.

352 4.3 Neural Networks

353 All NNs were trained in Python 3.6 using the Keras library (v2.3.1; TensorFlow v1.15.0 backend). Hardware
354 used was a Intel® Xeon® E2286G CPU and a NVIDIA GeForce RTX 2080. Training of all models took less
355 than 20 minutes.

356 Two general types of NNs were investigated: dense networks where there is a weighted connection
357 between every node in subsequent layers, and CNNs. The number of nodes in each hidden layer of a dense
358 network was defined as half of the previous layer. For example, with two hidden layers, hidden layer 1
359 would have 3,168 nodes and layer 2 would have 1,584.

360 A general schematic of the CNN structure is shown in Figure 6. Convolutional layers involved training
 361 a set of 2D filters or kernels, which are weighting functions multiplied with input values in a specific spatial
 362 window. Filters are smaller than the input dimensions and are slid across the entire input to produce
 363 feature mapping of the input. Since one filter or weighting function is used for the whole input space,
 364 fewer trainable parameters are needed than dense networks. It also gives rise to feature invariance since
 365 the trained filter can identify a feature in any position of the input space. Max pooling layers look for the
 366 maximum value within a spatial window and then uses that maximum value to represent the output over
 367 that spatial window, effectively reducing dimensionality. For more details on the mathematics of CNNs
 368 and the training process, see LeCun et al. (1998)⁴⁵ and Goodfellow et al., (2016).²⁸ CNN layers were
 369 considered as a set of convolutional layers followed by max pooling. Following convolutional layers, the
 370 structure was flattened, where pooled feature maps are vectorized into a dense hidden layer followed by
 371 a single output node. A varied number of structures were investigated to identify changes in performance
 372 based on depth (number of layers), size of convolutional filters (spatial window), number of max pooling
 373 layers, and number of convolutional filters. A summary of these structures is presented in Table 2.
 374



375 **Figure 6** General schematic of convolutional neural network structure. The number of convolutional layers
 376 as well as the number of layers (convolution + max pooling) can be varied.
 377

378 **Table 2** Descriptions of the general model types used.
 379

Model description	Num. hidden layers	Hidden layer	Output layers
Dense	n	Dense ($nodes = 0.5 \cdot previous\ layer$) -> Batch Normalization -> Activation (elu)	Output
CNN 1 layer, no max pooling	1	Convolution 2D ($number\ of\ filters, size\ of\ filter\ (x, x)$) -> Batch Normalization -> Activation (elu)	Flatten -> Output
CNN n layers	n	Convolution 2D ($number\ of\ filters, size\ of\ filter\ (x, x)$) -> Batch Normalization -> Activation (elu) -> Max Pooling 2D (2,2)	Flatten -> Output
CNN n convolutions, 4 layers	4	$n \times \{$ Convolution 2D ($number\ of\ filters, size\ of\ filter\ (x, x)$) -> Batch Normalization -> Activation (elu) $\} ->$ Max Pooling 2D (2,2)	Flatten -> Output

381

382 Common between all structures and types of NNs was the use of batch normalization to speed up
 383 training,⁴⁶ followed by activation using an exponential linear unit (elu) activation function (equation 1).
 384

$$385 \quad ELU: f(x) = \begin{cases} x, & x > 0 \\ e^x - 1, & x \leq 0 \end{cases} \quad (\text{equation 1})$$

386
 387 All networks were trained with a mean squared error loss function coupled with $L2$ regularization to
 388 prevent overfitting (equation 2). The Adam optimization algorithm was used for all training.
 389

$$390 \quad Loss = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda \|w^2\| \quad (\text{equation 2})$$

391 Where, y_i is the network output for sample i (prediction)
 392 \hat{y}_i is the true value for sample i
 393 N is the total number of samples
 394 λ is a hyperparameter to control $L2$ regularization (set to 0.01)
 395 w are all the network weights
 396

397 Prediction accuracy was determined on a test set (20% of all data, $n=28$) that was not used for training
 398 the network. The metric used to assess predictive performance was mean absolute error (MAE), primarily
 399 since it provides a metric in the same units used in analysis and is more easily interpreted.

400 4.4 Occlusion method

401 An occlusion approach was used to generate heat maps of spectral areas that most influence prediction
 402 accuracy. After training a network, test data was modified by iteratively setting a spectral area or patch
 403 equal to 0. These occluded or corrupted training samples are then fed through the network to produce a
 404 prediction of DBP concentration (Table 3).
 405

406 **Table 3** Description of the occlusion method used to identify spectral heat maps or areas of importance.

Occlusion method	
1:	Train a network using original training data (X_{train})
2:	Predict outputs (y_{test}) using original test data (X_{test})
3:	For $t = 1$ to 20,000 do
4:	Randomly select EEM patch from X_{test}
5:	Set patch = 0 to create corrupted test set, $X_{occluded}$
6:	Predict outputs ($y_{occluded}$) using $X_{occluded}$
7:	Average error calculated over all test data. $Error = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} (y_{test} - y_{occluded})$, where N_{test} is the number of samples in the test set.
8:	Check all excitation/emission were included in random selection
9:	Calculate average error over all iterations for each excitation/emission pair

407 The difference between initial predictions and occluded predictions provided an estimate of the
 408 importance of the occluded patch. If initial predictions and occluded predictions are identical or close, the
 409 trained network is not relying on that spectral area to estimate DBPs. On the other hand, if the error is
 410 high, the occluded region is influential on the accurate prediction of DBP levels.

411 References

- 412 1. Storey, M. V., van der Gaag, B. & Burns, B. P. Advances in on-line drinking water quality monitoring
413 and early warning systems. *Water Research* **45**, 741–747 (2011).
- 414 2. Faassen, S. & Hitzmann, B. Fluorescence Spectroscopy and Chemometric Modeling for Bioprocess
415 Monitoring. *Sensors* **15**, 10271–10291 (2015).
- 416 3. Beutel, S. & Henkel, S. In situ sensor techniques in modern bioprocess monitoring. *Applied*
417 *Microbiology and Biotechnology* vol. 91 1493–1505 (2011).
- 418 4. Bridgeman, J., Bieroza, M. & Baker, A. The application of fluorescence spectroscopy to organic matter
419 characterisation in drinking water treatment. *Reviews in Environmental Science and Biotechnology*
420 vol. 10 277–290 (2011).
- 421 5. Murphy, K. R., Stedmon, C. A. & Bro, R. Chemometric analysis of organic matter fluorescence. in
422 *Aquatic organic matter fluorescence* 339–375 (2014). doi:10.13140/2.1.2595.8080.
- 423 6. Wagner, E. D. & Plewa, M. J. CHO cell cytotoxicity and genotoxicity analyses of disinfection by-
424 products: An updated review. *Journal of Environmental Sciences* **58**, 64–76 (2017).
- 425 7. Guilherme, S., Dorea, C. C. & Rodriguez, M. J. Decision-making scheme for disinfection by-product
426 monitoring intended for small drinking water systems. *Environ. Sci.: Water Res. Technol.* **3**, 366–376
427 (2017).
- 428 8. Chen, B. & Westerhoff, P. Predicting disinfection by-product formation potential in water. *Water*
429 *Research* **44**, 3755–3762 (2010).
- 430 9. Kulkarni, P. & Chellam, S. Disinfection by-product formation following chlorination of drinking water:
431 Artificial neural network models and changes in speciation with treatment. *Science of The Total*
432 *Environment* **408**, 4202–4210 (2010).
- 433 10. Lin, H. *et al.* Radial basis function artificial neural network able to accurately predict disinfection by-
434 product levels in tap water: Taking haloacetic acids as a case study. *Chemosphere* **248**, 125999 (2020).
- 435 11. Singh, K. P. & Gupta, S. Artificial intelligence based modeling for predicting the disinfection by-
436 products in water. *Chemometrics and Intelligent Laboratory Systems* **114**, 122–131 (2012).
- 437 12. Sadiq, R. & Rodriguez, M. J. Disinfection by-products (DBPs) in drinking water and predictive models
438 for their occurrence: a review. *Science of The Total Environment* **321**, 21–46 (2004).
- 439 13. Matilainen, A. *et al.* An overview of the methods used in the characterisation of natural organic matter
440 (NOM) in relation to drinking water treatment. *Chemosphere* **83**, 1431–1442 (2011).
- 441 14. Murphy, K. R., Bro, R. & Stedmon, C. A. Chemometric Analysis of Organic Matter Fluorescence. in
442 *Aquatic Organic Matter Fluorescence* (eds. Coble, P., Lead, J., Baker, A., Reynolds, D. M. & Spencer, R.
443 G. M.) 339–375 (Cambridge University Press, 2014). doi:10.1017/CBO9781139045452.016.
- 444 15. Trueman, B. F., MacIsaac, S. A., Stoddart, A. K. & Gagnon, G. A. Prediction of disinfection by-product
445 formation in drinking water via fluorescence spectroscopy. *Environ. Sci.: Water Res. Technol.* **2**, 383–
446 389 (2016).
- 447 16. Pifer, A. D. & Fairey, J. L. Improving on SUVA₂₅₄ using fluorescence-PARAFAC analysis and asymmetric
448 flow-field flow fractionation for assessing disinfection byproduct formation and control. *Water*
449 *Research* **46**, 2927–2936 (2012).
- 450 17. Peleato, N. M., Legge, R. L. & Andrews, R. C. Neural networks for dimensionality reduction of
451 fluorescence spectra and prediction of drinking water disinfection by-products. *Water Research* **136**,
452 84–94 (2018).
- 453 18. Xu, X. *et al.* EEM-PARAFAC characterization of dissolved organic matter and its relationship with
454 disinfection by-products formation potential in drinking water sources of northeastern China. *Science*
455 *of The Total Environment* **774**, 145297 (2021).

- 456 19. Ma, C., Xu, H., Zhang, L., Pei, H. & Jin, Y. Use of fluorescence excitation–emission matrices coupled
457 with parallel factor analysis to monitor C- and N-DBPs formation in drinking water recovered from
458 cyanobacteria-laden sludge dewatering. *Science of The Total Environment* **640–641**, 609–618 (2018).
- 459 20. Yang, X., Shang, C., Lee, W., Westerhoff, P. & Fan, C. Correlations between organic matter properties
460 and DBP formation during chloramination. *Water Research* **42**, 2329–2339 (2008).
- 461 21. Peiris, R. H. *et al.* Identifying fouling events in a membrane-based drinking water treatment process
462 using principal component analysis of fluorescence excitation-emission matrices. *Water Research* **44**,
463 185–194 (2010).
- 464 22. Shutova, Y., Baker, A., Bridgeman, J. & Henderson, R. K. Spectroscopic characterisation of dissolved
465 organic matter changes in drinking water treatment: From PARAFAC analysis to online monitoring
466 wavelengths. *Water Research* **54**, 159–169 (2014).
- 467 23. Murphy, K. R., Stedmon, C. A., Graeber, D. & Bro, R. Fluorescence spectroscopy and multi-way
468 techniques. PARAFAC. *Analytical Methods* **5**, 6557–6566 (2013).
- 469 24. Sharpless, C. M. & Blough, N. V. The importance of charge-transfer interactions in determining
470 chromophoric dissolved organic matter (CDOM) optical and photochemical properties. *Environ. Sci.:
471 Processes Impacts* **16**, 654–671 (2014).
- 472 25. Kothawala, D. N., Murphy, K. R., Stedmon, C. A., Weyhenmeyer, G. A. & Tranvik, L. J. Inner filter
473 correction of dissolved organic matter fluorescence. *Limnology and Oceanography: Methods* **11**, 616–
474 630 (2013).
- 475 26. Bro, R. PARAFAC. Tutorial and applications. in *Chemometrics and Intelligent Laboratory Systems* vol.
476 38 149–171 (1997).
- 477 27. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
- 478 28. Goodfellow, I., Bengio, Y. & Courville, A. *Deep Learning*. (MIT Press, 2016).
- 479 29. Zhang, X. *et al.* Understanding the learning mechanism of convolutional neural networks in spectral
480 analysis. *Analytica Chimica Acta* **1119**, 41–51 (2020).
- 481 30. Bjerrum, E. J., Glahder, M. & Skov, T. Data Augmentation of Spectral Data for Convolutional Neural
482 Network (CNN) Based Deep Chemometrics. *arXiv:1710.01927 [cs]* (2017).
- 483 31. Malek, S., Melgani, F. & Bazi, Y. One-dimensional convolutional neural networks for spectroscopic
484 signal regression. *Journal of Chemometrics* **32**, e2977 (2018).
- 485 32. Chatzidakis, M. & Botton, G. A. Towards calibration-invariant spectroscopy using deep learning. *Sci
486 Rep* **9**, 2126 (2019).
- 487 33. Simonyan, K. & Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition.
488 *arXiv:1409.1556 [cs]* (2015).
- 489 34. Chen, W., Westerhoff, P., Leenheer, J. A. & Booksh, K. Fluorescence Excitation–Emission Matrix
490 Regional Integration to Quantify Spectra for Dissolved Organic Matter. *Environ. Sci. Technol.* **37**,
491 5701–5710 (2003).
- 492 35. Roccaro, P., Vagliasindi, F. G. A. & Korshin, G. V. Changes in NOM Fluorescence Caused by Chlorination
493 and their Associations with Disinfection by-Products Formation. *Environ. Sci. Technol.* **43**, 724–729
494 (2009).
- 495 36. Lakowicz, J. R. *Principles of Fluorescence Spectroscopy*. (Springer Science & Business Media, 2013).
- 496 37. Dickenson, E. R. V., Summers, R. S., Croué, J.-P. & Gallard, H. Haloacetic acid and Trihalomethane
497 Formation from the Chlorination and Bromination of Aliphatic β -Dicarbonyl Acid Model Compounds.
498 *Environ. Sci. Technol.* **42**, 3226–3233 (2008).
- 499 38. Zeng, T. & Arnold, W. A. Clustering Chlorine Reactivity of Haloacetic Acid Precursors in Inland Lakes.
500 *Environ. Sci. Technol.* **48**, 139–148 (2014).
- 501 39. Lavonen, E. E. *et al.* Tracking changes in the optical properties and molecular composition of dissolved
502 organic matter during drinking water production. *Water Research* **85**, 286–294 (2015).

- 503 40. Hong, H. C., Wong, M. H. & Liang, Y. Amino Acids as Precursors of Trihalomethane and Haloacetic Acid
504 Formation During Chlorination. *Arch Environ Contam Toxicol* **56**, 638–645 (2009).
- 505 41. Hua, L.-C., Lin, J.-L., Chen, P.-C. & Huang, C. Chemical structures of extra- and intra-cellular algogenic
506 organic matters as precursors to the formation of carbonaceous disinfection byproducts. *Chemical*
507 *Engineering Journal* **328**, 1022–1030 (2017).
- 508 42. Nemani, V. A., Taylor-Edmonds, L., Peleato, N. M. & Andrews, R. C. Impact of operational parameters
509 on biofiltration performance: organic carbon removal and effluent turbidity. *Water Science &*
510 *Technology: Water Supply* **16**, 1683–1692 (2016).
- 511 43. Peleato, N. M., Sidhu, B. S., Legge, R. L. & Andrews, R. C. Investigation of ozone and peroxone impacts
512 on natural organic matter character and biofiltration performance using fluorescence spectroscopy.
513 *Chemosphere* **172**, 225–233 (2017).
- 514 44. Lawaetz, A. J. & Stedmon, C. A. Fluorescence Intensity Calibration Using the Raman Scatter Peak of
515 Water. *Applied Spectroscopy* **63**, 936–940 (2009).
- 516 45. Lecun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-based learning applied to document
517 recognition. *Proceedings of the IEEE* **86**, 2278–2324 (1998).
- 518 46. Ioffe, S. & Szegedy, C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal
519 Covariate Shift. in *International Conference on Machine Learning* 448–456 (PMLR, 2015).
- 520

521 Acknowledgements

522 This work was funded in part by the Canadian Water Network and the Natural Sciences and
523 Engineering Research Council of Canada (NSERC).

524 Author Contributions

525 NP designed the work, wrote the scripts, conducted the analysis, and wrote the manuscript