

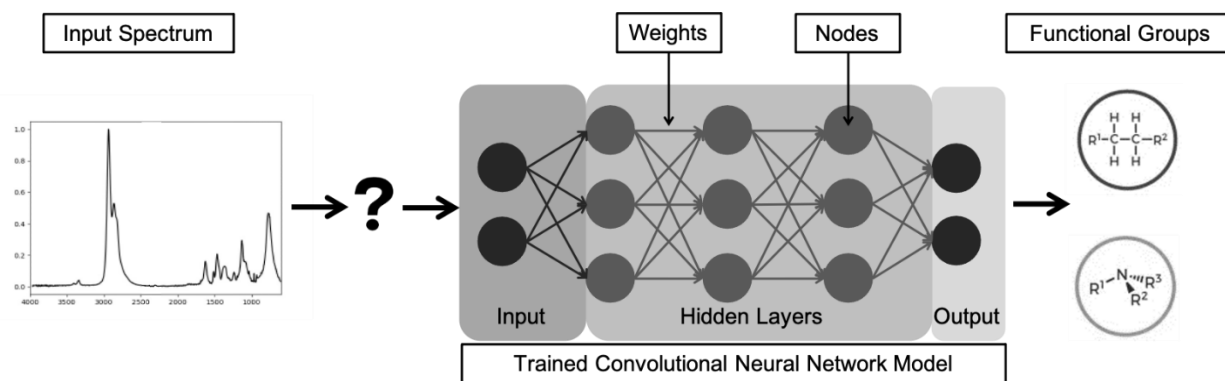
## **Functional group identification for FTIR spectra using image-based machine learning models**

Abigail A. Enders, Nicole M. North, Chase M. Fensore, Juan E. Velez-Alvarez, Heather C. Allen\*

Department of Chemistry & Biochemistry, The Ohio State University, Columbus, Ohio 43210, United States

### **Abstract**

Fourier Transform Infrared Spectroscopy (FTIR) is a ubiquitous spectroscopic technique. Spectral interpretation is a time-consuming process, but it yields important information about functional groups present in compounds and in complex substances. We develop a generalizable model via a machine learning (ML) algorithm using Convolutional Neural Networks (CNNs) to identify the presence of functional groups in gas phase FTIR spectra. The ML models will reduce the amount of time required to analyze functional groups and facilitate interpretation of FTIR spectra. Through web scraping, we acquire intensity-frequency data from 8728 gas phase organic molecules within the NIST spectral database and transform the data into images. We successfully train models for 15 of the most common organic functional groups, which we then determine via identification from previously untrained spectra. These models serve to expand the application of FTIR measurements for facile analysis of organic samples. Our approach was done such that we have broad functional group models that inference in tandem to provide full interpretation of a spectrum. We present the first implementation of ML using image-based CNNs for predicting functional groups from a spectroscopic method.



## Introduction

The anthropogenic impact on the climate and environment has prompted the analysis and detection of pollutants or contaminants with Fourier Transform Infrared Spectroscopy (FTIR), such as microplastics in waters<sup>1,2</sup> and table salts<sup>3</sup>, nitrates from agricultural fertilizers in soil<sup>4-6</sup>, and polyaromatic hydrocarbons in the ocean's surface<sup>7,8</sup>. The diversity of the chemical composition of the pollutants and the central fundamental technique of FTIR underscores the importance of a computational method for improved throughput of spectral analysis. The bottleneck is most frequently the assignment of peaks to relevant functional groups.<sup>9,10</sup>

Functional groups describe and define the physical and chemical properties of compounds.<sup>11,12</sup> Identification of many organic groups is accomplished via FTIR due to the associated unique vibrational frequencies.<sup>13,14</sup> Large numbers of spectra are time consuming to analyze and require expert chemist analysis to determine present composition. This limits the application of FTIR spectral techniques as a sampling method for functional group elucidation. There is thus an unexplored, yet applicable field of FTIR spectra interpretation through statistical methods. Progress towards machine learning (ML) methods for environmental pollutant analysis has been explored for specific, targeted applications.<sup>9,15-17</sup> Generalizable functional group ML models would increase the utility of FTIR sample screening in environmental and other chemistry applications.<sup>18,19</sup>

In this study, we investigate the implementation of convolutional neural networks (CNNs)<sup>20</sup> to identify functional groups present in FTIR spectra. By limiting spectral preprocessing, we explore a minimalistic approach to allow the network to learn spectral patterns for successful recognition of the fifteen most common organic functional groups (Table 1).

**Table 1.** Functional groups for which successful models were trained.

alcohol	alkane	alkene	alkyne	amide
ether	acyl halide	alkyl halide	methyl	ketone
carboxylic acid	nitrile	amine	aromatic	ester

Machine learning (ML) serves to address a need for quick identification of spectral components.<sup>21</sup> To date, the use of a CNN to broadly classify functional groups has not been reported. CNNs work by having layers of nodes called neurons, these neurons can be trained on data to identify spectral components that were observed in the training data in new spectra. The algorithm works to minimize a loss function; this is done by comparing answers given by the CNN to the true answers from a training data set. The difference between the reported and the true presence of a group constitute the loss function. The training data set is a randomly segmented subset of spectra that the CNN uses to learn and adjust neuron weights.

CNNs expand upon artificial neural networks (ANNs) by using mathematical convolutions to provide convolved data to the following neuron. Each neuron has a receptive field for which it convolves the information, similar to how a human brain has regions of neurons designated for processing specific

information. CNNs significantly reduce the number of neurons per pixel that a traditional feed-forward network requires to capture the complexity of an image. Thus, CNNs are a sophisticated solution to the alternative complex network required to machine learn images by capturing the spatial and temporal uniqueness of images.

We probe the effectiveness of image recognition ML as a facile solution to FTIR spectra interpretation. The information contained in a spectrum is most often presented to chemist as a 2D image, therefore it is desirable to develop models that learn via similar spectral visualization.<sup>22</sup> Previous implementations of FTIR ML for functional group identification have limited,<sup>23</sup> averaged,<sup>24</sup> and segmented<sup>23,25</sup> spectral data to reduce information used during training. The computational resources available today make this an unnecessary and limiting feature. We include all available spectral data from 4000 to 600  $\text{cm}^{-1}$  to reduce any biases on the learning process.

Current methods for spectral processing and interpretation are limited to library searching software<sup>26</sup> and highly specific questions using implementations of ML including: Support Vector Machines,<sup>9,27</sup> k-Nearest Neighbors,<sup>28,29</sup> and Principal Component Analysis (PCA)<sup>9,27,28</sup> or Factor Analysis<sup>30</sup>. Library searching methods require a pre-existing and transferrable database for searching spectra. The initial creation of libraries requires an intensive endeavor for collecting a large enough spectral repository. Once implemented, libraries cannot extrapolate beyond those included in the software. The size of libraries is not of significant concern for storage, but it is a cumbersome feature for application compatibility and relative use-to-memory consumption. ML does not require transfer of training data to the user and can predict beyond the data used for training.

The use of ML to resolve challenging implementations of FTIR spectra (e.g., extremely large datasets, continuous analysis) has become of interest as increased processing power makes it possible to train and infer (interpret an unknown spectra) with complex algorithms.<sup>31-34</sup> However, these highly specific models are only applicable in the setting in which they are developed because the training is completed

on a narrow range of examples. To increase the amount of available training spectra or improve further calculations, ML algorithms in tandem with molecular dynamics have been explored.<sup>34–36</sup>

Previous applications<sup>37–39</sup> of ML have employed data preprocessing prior to training with unsupervised ML methods, such as PCA<sup>27</sup>, which reduces the information in the training data. Spectral preprocessing is an unnecessary strategy with the advances in ML and doing so would limit the transferability of the final model to broader applications. Selecting spectral regions of interest can lead to a reduction of learning to memorization by the computer, meaning that rather than making a general model that can make inferences on novel spectra the model overfits the training data, meaning that it does well on what it has seen before but performs poorly on new data. Showing select data based on human evaluation increases the time required by an expert and potentially introduces overfitting. These unintended consequences include not allowing the computer to learn important spectral features by evaluating the entire spectrum rather than regions. While there are regions of relative disinterest to the chemist, it is not sufficient to ignore them in training. The absence of a peak is equally as informative as the presence of another. A recent application of ML successfully implemented broader methods for functional group analysis, however the authors utilize a multilayer perceptron ML method with an autoencoder and train using two sources of data: FTIR and MS spectra.<sup>40</sup>

In our work, we create separate functional group models that are run simultaneously resulting in complete analysis of FTIR spectra. The use of individual functional group models presents a robust approach to establish a broad but precise computational analysis of spectra. Training a model for each functional group improves the overall accuracy attainable because each model is focused on a binary question: is this functional group present? The training of individual models does not impede speed of spectrum analysis achieved and results are provided succinctly. By approaching the classification of spectra via the proposed method, we reduce the likelihood that the model learns a connection between functional groups that is not chemically relevant. In other words, one present functional group does

not indicate another group's presence or absence. Individually trained models reduce the potential for this and improve the overall accuracy by posing a simplified question. Here we develop effective and accurate FTIR ML models that apply to broader questions, limit spectral preprocessing, and provide the entire spectrum to the algorithm.

## Methods

*Python scripts.* All Python scripts can be accessed from our repository at this address: <https://github.com/Ohio-State-Allen-Lab/FTIRMachineLearning>. The FTIR spectra are property of NIST and can be accessed through their website. The implementation of Inception V3<sup>41</sup> is modified for our use and the original source is linked on our repository with the published modified version. The computational procedure is described in detail in the SI and is documented in each Python script.

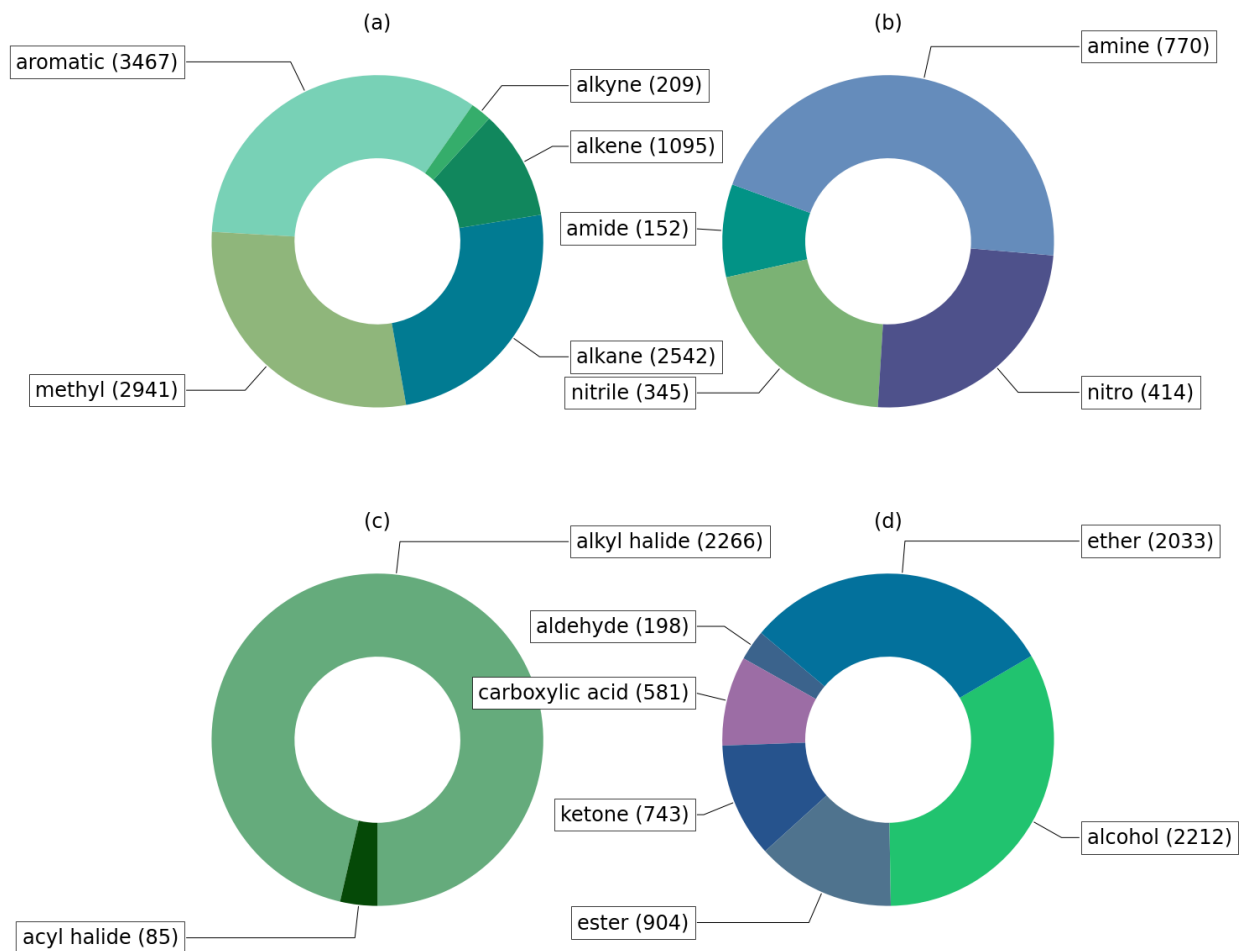
*Spectra collection.* Data was obtained from the National Institute for Science and Technology Chemistry WebBook via a web scraping implementation in Selenium using the CAS number identifier from the official list of compounds in the WebBook.<sup>42</sup> When a compound had an FTIR spectrum, the file, in jcamp-dx format, was retrieved and stored with the CAS number as the filename. A total of 8,728 spectra from pure compounds in gas phase were obtained. Each spectrum's InChI key was saved in a collective text file.

*Data pre-processing.* Only spectra in absorbance and wavenumbers were used for training models. Each spectrum was evaluated to ensure it was in absorbance and wavenumbers via a Python script. Files in transmission or wavelength were relocated to a distinct directory to preserve all spectra obtained from web scraping. Files in the correct mode were converted from jcamp-dx to csv. Once converted, each spectrum was normalized so that the maximum peak height was 1. Normalized spectra were saved as jpg images.

*Labeling.* Functional groups were identified via the InChI key. Using SMARTS functional group identifiers, each spectrum's key was parsed to return binary indicators. Present functional groups are

labelled as “1” and absent as “0”. Results were saved in one spreadsheet with CAS number’s as spectrum and file identifiers. Spectra were copied into directories based on presence or absence of a functional group. This method allows one compound with multiple functional groups present to be copied into the directory for each group. Each of the 17 functional groups had two directories: positive and negative cases. Positive cases include the functional group and negative cases do not contain the group. Randomly ten photos, five from positive and negative, for each functional group were reserved for validation. Then, the directory containing more instances for a given group was reduced randomly until both directories contained the same number of spectra.

*Machine learning.* A convolutional neural network (CNN) for image recognition was employed. A unique model was trained for each functional group containing two classes. The functional groups and the number of images in the positive cases are presented in the SI. The architecture, Inception V3, was accessed from the available models on the Google TensorFlow library. Each model was trained for 10,000 steps at a learning rate of 0.01, using an initialized version of Inception V3 and training the last layer of the model graph. The initialized parameters reduce the time and computational power required to train a custom model. It took five hours to train the fifteen models and classification of an unknown spectrum requires one minute.



**Figure 1.** Number of spectra used to train each functional group model, (a) carbon-containing, (b) nitrogen-containing, (c) halide-containing, and (d) oxygen-containing. The number of images is equivalent for the positive and negative cases used in training and testing.

*Accuracy and loss.* Accuracy and cross entropy (loss) for both training and test models was obtained and saved as csv files. The final accuracies and entropies for training and test results from each model are investigated to identify any anomalies.

*Classification of Validation Data.* When spectra were classified, the models were all called upon to infer (determine the functional groups present) and a final result was provided. The ten reserved spectra



were analyzed via the respective models they were withheld from to examine the learning quality of the algorithm. Confusion matrices<sup>43</sup> were used to represent the true and predicted functional group for the seventeen models.

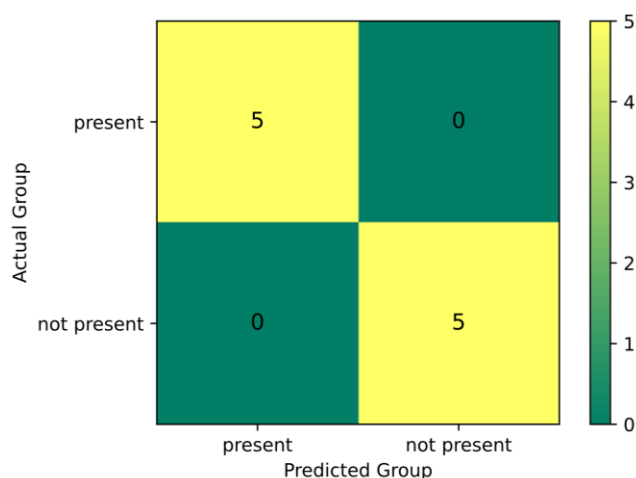
## Results and Discussion

Out of 17 functional group models trained, 15 effectively identify functional groups (Table 2). Accuracy and cross entropy results from the last step of training are reported for the train and validation process. The two functional group models that underperform are aldehyde and nitrile, based on model prediction of untrained spectra (as seen in the SI). We define underperforming as misidentifying more than 60% of test cases. The training accuracy is a measure of how well the model classifies the training data, which it used to train the network. A higher training accuracy indicates that the model is learning the training spectra. Validation accuracy expresses the ability of the model to generalize to untrained spectra, which is determined by the number of correctly classified validation spectra. Thus, it is more meaningful to have a higher validation accuracy, albeit not a requirement for a successful inferencing model. Cross entropy is the loss function used to evaluate the final model and is defined as the logarithm of the likelihood of a correct assignment. Smaller cross entropy values indicate a model is well trained. We observe cross entropy for training is less than validation. Models are more likely to correctly inference spectra that have been used to train and adjust weights, in comparison to the validation spectra.

**Table 2.** Final accuracy and cross entropy for train and validation of each functional group model is presented in order of increasing number of training images.

	Accuracy		Cross Entropy	
	Train (%)	Validation (%)	Train	Validation
acyl halide	100	98	0.025347	0.143665
amide	100	70	0.060037	0.900095
aldehyde	100	80	0.04735	0.385665
alkyne	99	80	0.079596	0.332745
nitrile	97	65	0.17019	0.668148
nitro	98	89	0.12624	0.668148
carboxylic acid	98	98	0.070173	0.076216
ketone	93	76	0.228837	0.501178
amine	93	80	0.24136	0.494815
ester	97	83	0.111057	0.323917
alkene	85	68	0.407803	0.743543
ether	89	81	0.27397	0.443644
alcohol	90	86	0.236544	0.330591
alkyl halide	85	73	0.350733	0.531302
alkane	85	90	0.327755	0.265718
methyl	81	84	0.384048	0.358021
aromatic	92	89	0.199645	0.259044

A confusion matrix for each model was created by using spectra that have been withheld from training and testing data (as seen in the SI). A confusion matrix compares model assignments to the actual identities of the samples; it shows correct assignments along the trace of a matrix and false assignments off of the trace. Four models have perfect confusion matrices from classification of ten withheld images, five containing and not-containing functional group spectra examples. The presence or absence of carboxylic acid, aromatic, methyl, and ester functional groups are correctly identified in the withheld spectra (Figure 2).



**Figure 2.** Confusion matrix for carboxylic acid, aromatic, methyl, and ether functional group models.

The number of instances of each functional group occurring in the spectra varies significantly, with aromatic-containing spectra occurring most frequently with 3,467 images. In contrast, acyl halide has 85 spectra for training and testing the model. We explored the relationship between the number of images and the cross entropy and accuracy for training and testing results (Figure 3). Training accuracy decreases with increasing number of spectra (Table 3, SI). However, the final accuracy, determined by evaluating the unknown spectra for functional group identification, is not correlated to the number of images used for training. Our results indicate that the total number of training spectra does not affect the final performance of the models. The scattered, non-uniformity exhibited in Figure 3 (a) and (b) depict the deviation from a linear relationship between the number of spectra and accuracy and cross entropy for validation, confirming the number of images is not influencing the performance of the models. Training accuracy provides insight into how well the model has learned the training images for a functional group model. Counterintuitively, few spectra being trained for a functional group will result in a higher training accuracy because the model trains on the same spectra more frequently. The model memorizes or overfits functional groups resulting in a model incapable of extrapolating to new spectra.

However, from our results the challenge of limited training spectra does not implicate less accurate models. We confirmed this by investigating the relationship of number of images per class as a function of validation accuracy and cross entropy (Table 2). Models that have more spectra to train on have lower overall training accuracy but still perform well when analyzing unknown spectra. To investigate the linear correlation between the number of spectra used for training and the final accuracy and cross entropy, the Pearson's correlation coefficient is used (Figure 3). More linearly correlated relationships have a coefficient closer to one, where positive coefficients indicate a positive correlation and negative coefficients indicate a negative correlation. The coefficient for training accuracy and number of training spectra indicate they are indirectly correlated whereas the coefficient for training cross entropy and number of training spectra is positive, or positively correlated. However, the models with less training spectra show no correlation between final accuracy and ability to classify unknown spectra. Furthermore, both validation accuracy and cross entropy do not have significant linear correlation. While training results display correlation with the number of images, the validation data indicates that models are successful with a range of number of training spectra.

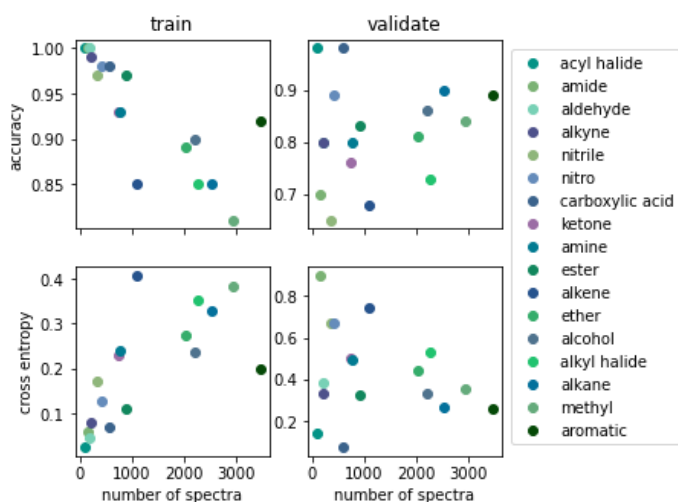
We can determine some of the underlying shortcomings in the model, from both spectroscopic and computational perspectives, by investigating two functional groups: aldehyde and carboxylic acid. The model results for aldehyde are promising for the training data but do not perform as effectively in validation and testing (Table 2, SI). The confusion matrix for carboxylic acid describes how well the model performs on spectra that have not been used for training or validation. We observe that the IR mode frequencies for the carboxylic acid and aldehyde affect the performance of the model, in addition to the number of spectral examples available for training and validation.

Aldehyde C-H stretching frequency ( $2830\text{-}2695\text{ cm}^{-1}$ ) is commonly overlapping in organic spectra with other C-H bonds because it is a weaker mode (Table 3). The carbonyl stretch is also frequently unresolved in compounds that contain multiple oxygen atoms. The C-H bending mode is often weak,

in addition to being in the fingerprint region, which is a challenge to interpret due to the complexity. With these stretching and bending modes considered, it is reasonable to anticipate that an aldehyde functional group is challenging for the model to identify in spectra. In comparison, carboxylic acid functional groups are always correctly identified in spectra by the model. The model for carboxylic acids is well trained. As observed by the validation accuracy and cross entropy (Table 2), the carboxylic acid model has a more robust transferability to spectra it has never observed. We confirm the effectiveness of the model with a correct assignment of unknown spectra. In total, there are more carboxylic acid training spectra, and the IR modes are better resolved, especially the strong COO-H stretching, in comparison to aldehydes (Table 3).

**Table 3.** Aldehyde and carboxylic acid IR stretching and bending mode frequencies.<sup>44</sup>

	Mode	Frequency (cm <sup>-1</sup> )	Appearance
<b>Aldehyde</b>	C-H stretch	2830-2695	Weak, medium
	C=O stretch	1740-1720	Medium, strong
	C-H bend	1390-1380	Weak, medium
<b>Carboxylic acid</b>	O-H stretch	3300-2500	Strong, broad
	C=O stretch	1760	Strong
	O-H bend	1440-1395	Medium



**Figure 3.** Final train and validation accuracy and cross entropy as a function of the number of spectra used to train each functional group. Pearson's correlation coefficients (PCC) are inset in the plots for final accuracy and cross entropy of training and validation as a function of the number of spectra. The coefficients closer to  $\pm 1$  indicate that the train accuracy and cross entropy are linearly correlated (negative is inversely correlated and positive is directly correlated) to the number of spectra used in training. Validation accuracy and cross entropy are not linearly correlated to the spectral examples used.

From our results, we observe that the models are more accurate for functional groups when there are more training spectra examples for the functional group and IR peaks are well resolved. Albeit this is an intuitive result for a trained spectroscopist with respect to accuracy correlating to peak resolution, yet there is no precedent using a machine learning approach.

### **Conclusion**

We present a novel method for FTIR spectral interpretation using CNNs and the NIST database. Fifteen functional group models successfully and effectively classify unknown spectra in a facile method for spectral submission to interpretation. We find that the image recognition features inherent in CNNs are transferrable to a chemical-identification application. From our observations, we can conclude that CNNs are effective at identifying spectral features for classification and generalizable models are achievable with ample spectral examples. In future work, optimization for functional group identification with less spectral examples should be investigated to improve accuracy.

## **AUTHOR INFORMATION**

### **Corresponding Author**

\* Heather C. Allen – Department of Chemistry & Biochemistry, The Ohio State University, Columbus, Ohio 43210, United States; orcid.org/0000-0003-3120-6784; Phone: +1- 614-292-4707; Email: allen@chemistry.ohio-state.edu; Fax: +1-614-292-1685

### **Author Contributions**

A.E. conducted experiments and writing. N.N. contributed code and editing. C.F. wrote Selenium web scraping script. J.E.V. intellectual contributions. H.A. intellectual contributions.

### **Notes**

The authors declare no competing financial interest.

## **ACKNOWLEDGMENT**

The authors acknowledge R. L. Enders for helpful discussions.

## **REFERENCES**

- (1) Yonkos, L. T.; Friedel, E. A.; Perez-Reyes, A. C.; Ghosal, S.; Arthur, C. D. Microplastics in Four Estuarine Rivers in the Chesapeake Bay, U.S.A. *Environ. Sci. Technol.* **2014**, *48* (24), 14195–14202. <https://doi.org/10.1021/es5036317>.
- (2) Song, Y. K.; Hong, S. H.; Jang, M.; Kang, J. H.; Kwon, O. Y.; Han, G. M.; Shim, W. J. Large Accumulation of Micro-Sized Synthetic Polymer Particles in the Sea Surface Microlayer. *Environ. Sci. Technol.* **2014**, *48* (16), 9014–9021. <https://doi.org/10.1021/es501757s>.
- (3) Lee, H. J.; Song, N. S.; Kim, J. S.; Kim, S. K. Variation and Uncertainty of Microplastics in Commercial Table Salts: Critical Review and Validation. *J. Hazard. Mater.* **2021**, *402* (August 2020), 123743. <https://doi.org/10.1016/j.jhazmat.2020.123743>.

- (4) Bogard, J. S.; Johnson, S. A.; Kumar, R.; Cunningham, P. T. Quantitative Analysis of Nitrate Ion in Ambient Aerosols by Fourier-Transform Infrared Spectroscopy. *Environ. Sci. Technol.* **1982**, *16* (3), 136–140. <https://doi.org/10.1021/es00097a004>.
- (5) Anil, I.; Golcuk, K.; Karaca, F. ATR-FTIR Spectroscopic Study of Functional Groups in Aerosols: The Contribution of a Saharan Dust Transport to Urban Atmosphere in Istanbul, Turkey. *Water, Air, Soil Pollut.* **2014**, *225* (3). <https://doi.org/10.1007/s11270-014-1898-9>.
- (6) Linker, R.; Shmulevich, I.; Kenny, A.; Shaviv, A. Soil Identification and Chemometrics for Direct Determination of Nitrate in Soils Using FTIR-ATR Mid-Infrared Spectroscopy. *Chemosphere* **2005**, *61* (5), 652–658. <https://doi.org/10.1016/j.chemosphere.2005.03.034>.
- (7) Hardy, J. T. The Sea Surface Microlayer: Biology, Chemistry and Anthropogenic Enrichment. *Progress in Oceanography*. 1982. [https://doi.org/10.1016/0079-6611\(82\)90001-5](https://doi.org/10.1016/0079-6611(82)90001-5).
- (8) Wurl, O.; Obbard, J. P. A Review of Pollutants in the Sea-Surface Microlayer (SML): A Unique Habitat for Marine Organisms. *Mar. Pollut. Bull.* **2004**, *48* (11–12), 1016–1030. <https://doi.org/10.1016/j.marpolbul.2004.03.016>.
- (9) P. M. Michel, A.; E. Morrison, A.; L. Preston, V.; T. Marx, C.; C. Colson, B.; K. White, H. Rapid Identification of Marine Plastic Debris via Spectroscopic Techniques and Machine Learning Classifiers. *Environ. Sci. & Technol.* **2020**, *54* (17), 10630–10637. <https://doi.org/10.1021/acs.est.0c02099>.
- (10) Rezania, S.; Park, J.; Md Din, M. F.; Mat Taib, S.; Talaiekhosani, A.; Kumar Yadav, K.; Kamyab, H. Microplastics Pollution in Different Aquatic Environments and Biota: A Review of Recent Studies. *Mar. Pollut. Bull.* **2018**, *133* (May), 191–208. <https://doi.org/10.1016/j.marpolbul.2018.05.022>.
- (11) Carey, F. A.; Giuliano, R. M. *Organic Chemistry*, 10th ed.; McGraw-Hill Education, 2016.



- (12) Carey, F. A.; Sundberg, R. J. *Advanced Organic Chemistry Part A: Structure and Mechanisms*, 5th ed.; Springer US, 2007.
- (13) Barnes, R.; Gore, R. Infrared Spectroscopy. *Anal. Chem.* **1949**, *21* (1), 7–12. <https://doi.org/10.1021/ac60025a003>.
- (14) Peck, R. L. Characterization of Organic Compounds. *Anal. Chem.* **1950**, *22* (1), 121–126. <https://doi.org/10.1021/ac60037a027>.
- (15) Kedzierski, M.; Falcou-Préfol, M.; Kerros, M. E.; Henry, M.; Pedrotti, M. L.; Bruzard, S. A Machine Learning Algorithm for High Throughput Identification of FTIR Spectra: Application on Microplastics Collected in the Mediterranean Sea. *Chemosphere* **2019**, *234*, 242–251. <https://doi.org/10.1016/j.chemosphere.2019.05.113>.
- (16) Galimberti, D. R.; Bougueroua, S.; Mahé, J.; Tommasini, M.; Rijs, A. M.; Gaigeot, M. P. Conformational Assignment of Gas Phase Peptides and Their H-Bonded Complexes Using Far-IR/THz: IR-UV Ion Dip Experiment, DFT-MD Spectroscopy, and Graph Theory for Mode Assignment. *Faraday Discuss.* **2019**, *217*, 67–97. <https://doi.org/10.1039/c8fd00211h>.
- (17) Bougueroua, S.; Spezia, R.; Pezzotti, S.; Vial, S.; Quessette, F.; Barth, D.; Gaigeot, M. P. Graph Theory for Automatic Structural Recognition in Molecular Dynamics Simulations. *J. Chem. Phys.* **2018**, *149* (18). <https://doi.org/10.1063/1.5045818>.
- (18) Coe, J. V.; Nystrom, S. V.; Chen, Z.; Li, R.; Verreault, D.; Hitchcock, C. L.; Martin, E. W.; Allen, H. C. Extracting Infrared Spectra of Protein Secondary Structures Using a Library of Protein Spectra and the Ramachandran Plot. *J. Phys. Chem. B* **2015**, *119* (41), 13079–13092. <https://doi.org/10.1021/acs.jpcc.5b08052>.
- (19) Geiger, A.; Cao, Z.; Song, Z.; Ulcickas, J.; Simpson, G. Chapter 18. Autonomous Science: Big Data Tools for Small Data Problems in Chemistry; 2020; pp 450–487. <https://doi.org/10.1039/9781839160233-00450>.

- (20) Szegedy, C.; Vanhoucke, V.; Ioffe, S.; Shlens, J.; Wojna, Z. Rethinking the Inception Architecture for Computer Vision. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition.*; 2016.
- (21) Nalla, R.; Pinge, R.; Narwaria, M.; Chaudhury, B. Priority Based Functional Group Identification of Organic Molecules Using Machine Learning. *ACM Int. Conf. Proceeding Ser.* **2018**, 201–209. <https://doi.org/10.1145/3152494.3152522>.
- (22) Kowalski, B. R.; Jurs, P. C.; Isenhour, T. L.; Reilley, C. N.; Isenhour, T. L.; Jurs, P. C. Computerized Learning Machines Applied to Chemical Problems Interpretation of Infrared Spectrometry Data. *Anal. Chem.* **1969**, *41* (14), 1945–1949. <https://doi.org/10.1021/ac50159a026>.
- (23) Fessenden, R. J.; Györgyi, L. Identifying Functional Groups in IR Spectra Using an Artificial Neural Network. *J. Chem. Soc. Perkin Trans. 2* **1991**, 1755.
- (24) van Est, Q. C.; Schoenmakers, P. J.; Smits, J. R. M.; Nijssen, W. P. M. Practical Implementation of Neural Networks for the Interpretation of Infrared Spectra. *Vib. Spectrosc.* **1993**, *4* (3), 263–272. [https://doi.org/10.1016/0924-2031\(93\)80001-V](https://doi.org/10.1016/0924-2031(93)80001-V).
- (25) Wu, W.; Massart, D. L. Artificial Neural Networks in Classification of NIR Spectral Data: Selection of the Input. *Chemom. Intell. Lab. Syst.* **1996**, *35* (1), 127–135. [https://doi.org/10.1016/S0169-7439\(96\)00034-2](https://doi.org/10.1016/S0169-7439(96)00034-2).
- (26) Averill, D. F.; Baird, K. C.; Hopkins, L. L.; Yerkes, M. J. Fourier Transform Infrared Spectroscopy without an FTIR Spectrometer: Library Searching and Concise Storage of Spectra. *J. Chem. Inf. Comput. Sci.* **1990**, *30* (2), 133–136. <https://doi.org/10.1021/ci00066a006>.
- (27) Hugo J. M. dos Santos, V.; Do Canto Bruzza, E.; E. de Lima, J.; V. Lourega, R.; F. Rodrigues, L. Discriminant Analysis and Cluster Analysis of Biodiesel Fuel Blends Based on Fourier Transform Infrared Spectroscopy (FTIR). *Energy & Fuels* **2016**, *30* (6), 4905–4915. <https://doi.org/10.1021/acs.energyfuels.6b00447>.

- (28) Judge, K.; W. Brown, C.; Hamel, L. Sensitivity of Infrared Spectra to Chemical Functional Groups. *Anal. Chem.* **2008**, *80* (11), 4186–4192. <https://doi.org/10.1021/ac8000429>.
- (29) Renner, G.; Schmidt, T. C.; Schram, J. A New Chemometric Approach for Automatic Identification of Microplastics from Environmental Compartments Based on FT-IR Spectroscopy. *Anal. Chem.* **2017**, *89* (22), 12045–12053. <https://doi.org/10.1021/acs.analchem.7b02472>.
- (30) Golz, E. K.; Vander Griend, D. A. Modeling Methylene Blue Aggregation in Acidic Solution to the Limits of Factor Analysis. *Anal. Chem.* **2013**, *85* (2), 1240–1246. <https://doi.org/10.1021/ac303271m>.
- (31) H. Agbaria, A.; Beck Rosen, G.; Lapidot, I.; H. Rich, D.; Huleihel, M.; Mordechai, S.; Salman, A.; Kapelushnik, J. Differential Diagnosis of the Etiologies of Bacterial and Viral Infections Using Infrared Microscopy of Peripheral Human Blood Samples and Multivariate Analysis. *Anal. Chem.* **2018**, *90* (13), 7888–7895. <https://doi.org/10.1021/acs.analchem.8b00017>.
- (32) Sharaha, U.; Rodriguez-Diaz, E.; Sagi, O.; Riesenber, K.; Lapidot, I.; Segal, Y.; Bigio, I. J.; Huleihel, M.; Salman, A. Detection of Extended-Spectrum  $\beta$ -Lactamase-Producing Escherichia Coli Using Infrared Microscopy and Machine-Learning Algorithms. *Anal. Chem.* **2019**, *91* (3), 2525–2530. <https://doi.org/10.1021/acs.analchem.8b05497>.
- (33) Kedzierski, M.; Falcou-Préfol, M.; Kerros, M. E.; Henry, M.; Pedrotti, M. L.; Bruzard, S. A Machine Learning Algorithm for High Throughput Identification of FTIR Spectra: Application on Microplastics Collected in the Mediterranean Sea. *Chemosphere* **2019**, *234*, 242–251. <https://doi.org/10.1016/j.chemosphere.2019.05.113>.
- (34) Kananenka, A. A.; Yao, K.; Corcelli, S. A.; Skinner, J. L. Machine Learning for Vibrational Spectroscopic Maps. *J. Chem. Theory Comput.* **2019**, *15* (12), 6850–6858. <https://doi.org/10.1021/acs.jctc.9b00698>.

- (35) Morawietz, T.; Urbina, A. S.; Wise, P. K.; Wu, X.; Lu, W.; Ben-Amotz, D.; Markland, T. E. Hiding in the Crowd: Spectral Signatures of Overcoordinated Hydrogen-Bond Environments. *J. Phys. Chem. Lett.* **2019**, *10* (20), 6067–6073. <https://doi.org/10.1021/acs.jpcclett.9b01781>.
- (36) Zhai, Y.; Caruso, A.; Gao, S.; Paesani, F. Active Learning of Many-Body Configuration Space: Application to the Cs<sup>+</sup>-Water MB-Nrg Potential Energy Function as a Case Study. *J. Chem. Phys.* **2020**, *152* (14). <https://doi.org/10.1063/5.0002162>.
- (37) Doblies, A.; Boll, B.; Fiedler, B. Prediction of Thermal Exposure and Mechanical Behavior of Epoxy Resin Using Artificial Neural Networks and Fourier Transform Infrared Spectroscopy. *Polymers (Basel)*. **2019**, *11* (2), 363. <https://doi.org/10.3390/polym11020363>.
- (38) Lasch, P.; Stämmler, M.; Zhang, M.; Baranska, M.; Bosch, A.; Majzner, K. FT-IR Hyperspectral Imaging and Artificial Neural Network Analysis for Identification of Pathogenic Bacteria. *Anal. Chem.* **2018**, *90* (15), 8896–8904. <https://doi.org/10.1021/acs.analchem.8b01024>.
- (39) Lasch, P.; Diem, M.; Hänsch, W.; Naumann, D. Artificial Neural Networks as Supervised Techniques for FT-IR Microspectroscopic Imaging. *J. Chemom.* **2006**, *20* (5), 209–220. <https://doi.org/10.1002/cem.993>.
- (40) Fine, J. A.; Rajasekar, A. A.; Jethava, K. P.; Chopra, G. Spectral Deep Learning for Prediction and Prospective Validation of Functional Groups. *Chem. Sci.* **2020**, *11* (18), 4618–4630. <https://doi.org/10.1039/c9sc06240h>.
- (41) Linstrom, P. J.; Mallard, W. G. *NIST Chemistry WebBook, NIST Standard Reference Database Number 69*.
- (42) Stehman, S. V. Selecting and Interpreting Measures of Thematic Classification Accuracy. *Remote Sens. Environment* **1997**, *62*, 77–89.
- (43) Wade, L.; Simek, J. *Organic Chemistry*, 9th ed.; Pearson, 2016.

