

# 1           **Harnessing Semi-Supervised Machine Learning to Automatically Predict**

## 2                           **Bioactivities of Per- and Polyfluoroalkyl Substances (PFASs)**

3  
4   *Hyuna Kwon<sup>a</sup>, Zulfikhar A. Ali<sup>b</sup>, Bryan M. Wong<sup>a,b,\*</sup>*

5  
6           a) Department of Chemical & Environmental Engineering, University of California-Riverside, Riverside,  
7           CA 92521, United States

8           b) Department of Physics & Astronomy, University of California-Riverside, Riverside, CA 92521, United  
9           States

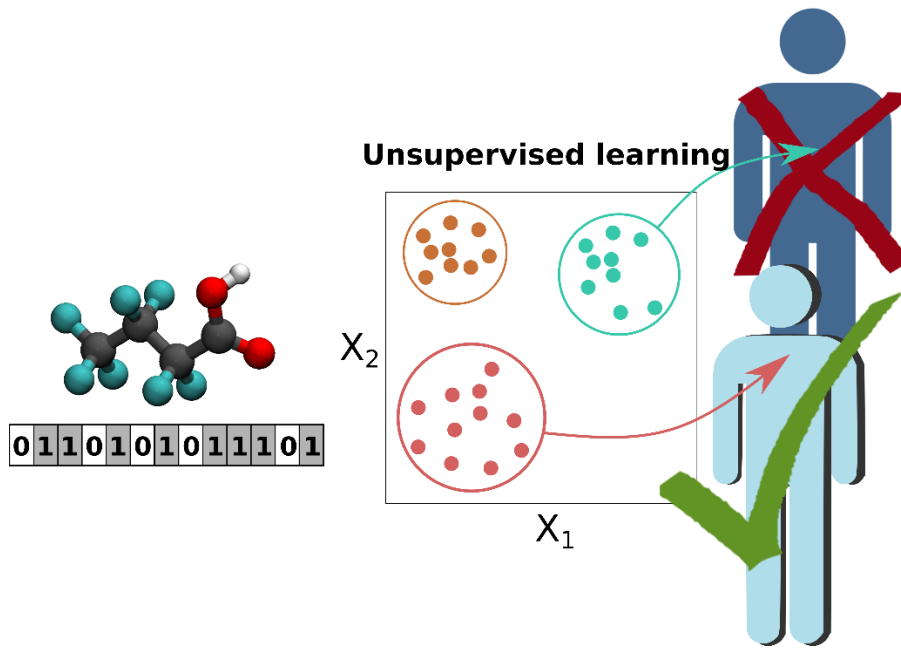
10  
11 \*Corresponding author. E-mail: [bryan.wong@ucr.edu](mailto:bryan.wong@ucr.edu); Web: <http://www.bmwong-group.com>

### 12 13 14 **Abstract**

15  
16 Many per- and polyfluoroalkyl substances (PFASs) pose significant health hazards due to their  
17 bioactive and persistent bioaccumulative properties. However, assessing the bioactivities of  
18 PFASs is both time-consuming and costly due to the sheer number and expense of *in vivo* and  
19 *in vitro* biological experiments. To this end, we harnessed new unsupervised/semi-supervised  
20 machine learning models to automatically predict bioactivities of PFAS in various human  
21 biological targets, including enzymes, genes, proteins, and cell lines. Our semi-supervised  
22 metric learning models were used to predict the bioactivity of PFASs found in the recent  
23 Organization of Economic Cooperation and Development (OECD) report list, which contains  
24 4,730 PFASs used in a broad range of industries and consumers. Our work provides the first  
25 semi-supervised machine learning study of structure-activity relationships for predicting  
26 possible bioactivities in a variety of PFAS species.

27  
28 **Keywords:** per- and polyfluoroalkyl substances, PFAS, machine learning, bioactivity, semi-  
29 supervised learning

30  
31 **Synopsis:** New machine learning techniques were used to automatically predict the  
32 bioactivities of PFAS in various human biological targets.



**Table of Contents Figure**

33  
34  
35

## 36 **Introduction**

37 Since the 1930s,<sup>1</sup> per- and polyfluoroalkyl substances (PFASs) have been used in several  
38 consumer products (including fire-fighting foams) due to their outstanding stability and  
39 water/oil repellent properties.<sup>2</sup> However, these compounds pose significant risks to the  
40 environment and biosystems. The presence of PFASs in surface water and groundwater can  
41 result in exposure to organisms, subsequently leading to accumulation in the body, with adverse  
42 effects on the liver, kidneys, blood, and immune system.<sup>2,3</sup> Because of these deleterious effects,  
43 there is a pressing need to identify and understand the bioactivity of PFAS-based compounds  
44 that can adversely affect human health.

45 For these reasons, several international groups including the Organization for Economic  
46 Cooperation and Development (OECD), United States Environmental Protection Agency,  
47 Food and Drug Administration, European Chemicals Agency, European Food Safety Authority,  
48 and Ministry of Ecology and Environment (China) continue to monitor PFASs that are  
49 produced in the global market.<sup>4,5</sup> According to a 2018 OECD report, more than 4,700 PFASs  
50 currently exist as manufacturers bring new forms of PFASs into industrial and consumer  
51 products (it is worth pointing out, however, that not all 4,700 structures exist in commerce).  
52 Nevertheless, among the wide varieties of PFAS molecules, the potential hazards of these new  
53 forms remain largely unknown.

54 Due to the sheer number of PFAS species, *in vivo* and *in vitro* biological experiments are  
55 both time-consuming and costly. As such, the construction of predictive and reliable  
56 quantitative-structure activity relationship (QSAR) models<sup>6-8</sup> is essential for assessing the  
57 bioactivities of these contaminants (even for PFAS species that are yet to be made). Specifically,  
58 a QSAR model that can accurately predict the bioactivities of PFASs can be harnessed to screen  
59 several of these contaminants, saving immense time and experimental resources. While there  
60 have been prior machine learning studies on PFAS molecules,<sup>9,10</sup> most of these approaches

61 used supervised learning techniques to suggest *general* structure-bioactivity trends after post-  
62 processing of the data (i.e., the focus was on aggregate data for all targets as opposed to  
63 analyzing chemical trends specific to each target).

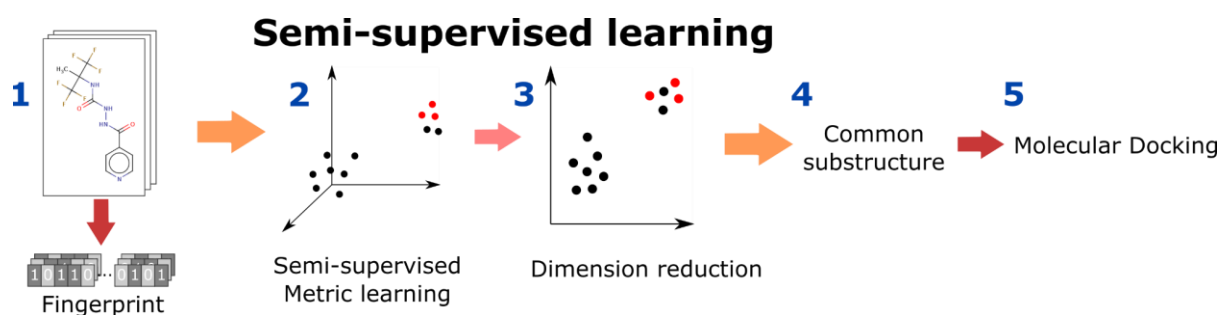
64 In this work, we present a new QSAR model using semi-supervised metric learning  
65 techniques to assess which functional groups affect bioactivities toward specific biological  
66 targets. Semi-supervised learning is a different machine learning approach that has the  
67 advantages of both supervised and unsupervised learning. It can be used on a dataset with  
68 primarily unlabeled data and only a few labeled data. Like unsupervised learning, it can also  
69 automatically cluster unlabeled data. Our approach is integrated with molecular docking  
70 calculations to predict possible bioactivities of PFAS molecules based on their chemical  
71 functional groups and specific biological targets (e.g., genes, proteins, or cell lines). Our  
72 approach first combines dimension reduction methods with clustering methods to classify  
73 PFASs based on their molecular structures. We then apply a semi-supervised metric learning  
74 method to improve classification accuracy. Finally, we use a molecular docking approach to  
75 shed light on the physicochemical reasons for their bioactivity. Our study provides the first  
76 unsupervised/semi-supervised learning approach for screening potentially bioactive PFAS  
77 molecules beyond conventional supervised learning or QSAR approaches.

78

## 79 **Methods**

80

81



82

83 **Figure 1:** Machine-learning-based workflow for QSAR construction to predict bioactivity of PFASs.

84

85 Our QSAR machine-learning framework, shown in Figure 1, utilizes four sequential steps  
86 followed by a reasoning/validation step: (1) collecting a training dataset from verified open-  
87 source databases, (2) encoding those compounds into molecular fingerprints, (3) clustering the  
88 data to predict chemical properties based on the molecular fingerprints and assessing the  
89 performance of the models, (4) evaluating the clustering by choosing the optimal model and  
90 predicting molecular groups responsible for bioactivity based on the clustering, and (5)  
91 molecular docking simulations to rationalize the role of the chemical functional groups. All of  
92 our machine learning algorithms are publicly available (see Supporting Information).

93 In our first step, we obtained datasets from comprehensive open-source databases,  
94 including PubChem's BioAssay,<sup>11</sup> Maximum Unbiased Validation,<sup>12</sup> Toxicology in the 21<sup>st</sup>  
95 Century,<sup>13</sup> beta-secretase 1,<sup>14</sup> and blood-brain barrier penetration datasets,<sup>15</sup> which are  
96 available from the Supporting Information of Ref. 10. We used two different datasets without  
97 further modification from Ref. 10: (1) the CF dataset, which includes substances containing at  
98 least one -CF- moiety (62,043 molecules), and (2) the C3F6 dataset, which includes  
99 substances containing a perfluoroalkyl moiety with three or more carbons (1,012 molecules).  
100 For both datasets, we used bioactivity data against 26 biological targets.

101 Encoding the compounds to molecular fingerprints followed next in our framework. We  
102 used the extended connectivity fingerprint (ECFP) featurization<sup>16</sup> with a default diameter of 4  
103 (i.e., ECFP4), which considers a maximum of four neighbors. ECFPs are topological molecular  
104 representations developed for substructure and similarity searching. By encoding molecular  
105 structures into fingerprints, we obtained a binary array with a constant length of 2,048, making  
106 it a convenient input for the unsupervised/semi-supervised learning models. Furthermore, since  
107 the simplified molecular-input line-entry system (SMILES) sequences for all PFAS molecules

108 are readily available, they can be easily converted into fingerprint-based representations using  
109 the RDKit software package.<sup>17</sup>

110 We then applied semi-supervised metric learning to the generated fingerprints by training  
111 machine learning models to predict the bioactivities of PFAS molecules by first (a) *reducing*  
112 *the dimension of the fingerprint datasets* and then (b) *classifying/clustering them* (see Figure  
113 1). Our QSAR model used a semi-supervised metric learning algorithm to automatically  
114 group/classify molecules with similar bioactivities. Metric learning has two main advantages:  
115 (1) its predictions are more efficient/accurate since the model distinctly separates new  
116 molecular representations according to their bioactivities (by reducing the distance metric  
117 between the same-labeled pair of data and increasing the distance between opposite-labeled  
118 pair of data), and (2) it automatically generates a vector-shaped representation from the  
119 molecular fingerprint and can be directly integrated with conventional dimension reduction  
120 methods. The final clusters were selected based on the best Silhouette score, which analyzes  
121 the distances of each data point to its cluster and neighboring clusters.<sup>18</sup> In short, a higher  
122 Silhouette score indicates more distinct and separated clusters. We then identified which  
123 substructures or molecular functional groups played essential roles in determining the  
124 bioactivity of the molecules.

125 Lastly, we conducted several molecular docking calculations using Autodock<sup>19</sup> to elucidate  
126 the physicochemical reasons for the bioactivity trends obtained from our QSAR model (i.e.,  
127 using ligand-protein binding conformations to rationalize the role of chemical substructures  
128 that induces bioactivity on biological targets.)

129

## 130 **Results and Discussion**

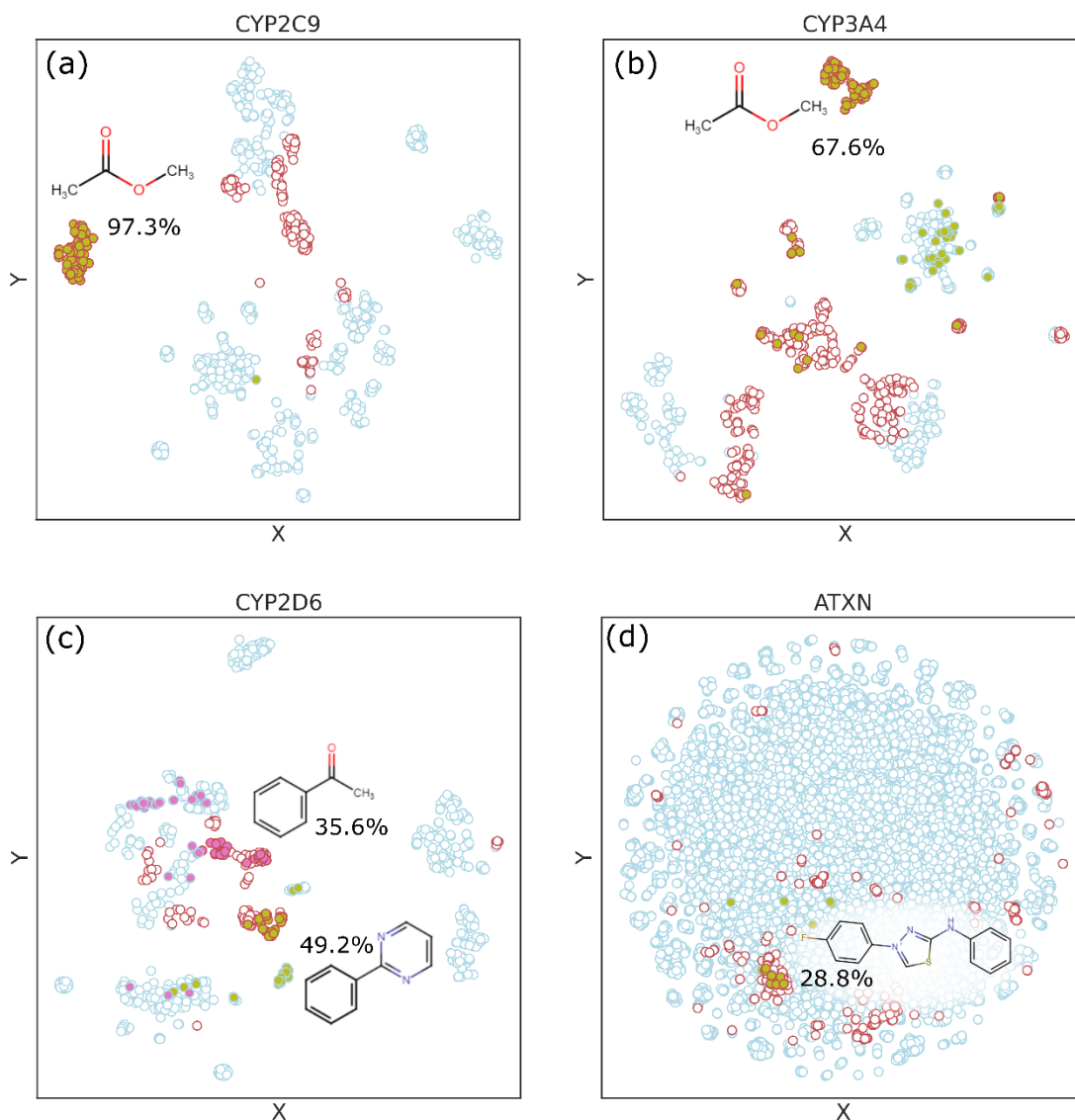
### 131 **3-1. Unsupervised vs. semi-supervised machine learning**

132 To systematically evaluate the performance of our semi-supervised metric approach, we  
133 first performed traditional unsupervised machine learning and compared the performance of  
134 the two models. To maintain a concise discussion of our results, the Supporting Information  
135 contains a detailed analysis and comparison of our unsupervised vs. semi-supervised machine  
136 learning results. Figure S1 shows our clustering results using unsupervised machine learning  
137 on the C3F6 dataset, and Figure S2 shows a comparison between the unsupervised and semi-  
138 supervised results using the CF dataset on two different targets. Table S3 summarizes the  
139 substructures that induce bioactivity as predicted from our unsupervised learning calculations.  
140 In summary, our extensive analyses in the Supporting Information showed that semi-  
141 supervised metric learning performed significantly better than unsupervised machine learning;  
142 as such, we only focus on the results of the former in this manuscript.

143  
144

### 145 **3-2. Semi-supervised metric learning**

146 Figure 2 displays true-positive ratios and classifications between bioactive/inactive  
147 molecules on four representative targets that show the best performance in the CF dataset using  
148 semi-supervised metric learning (for example, in Fig. 2a, we obtain a true-positive ratio of 97.3%  
149 by computing  $\frac{\text{number of molecules containing esters and are also bioactive}}{\text{number of ester-containing molecules in the cluster}}$ ). Using the Maximum  
150 Common Structure (MCS) module in the RDKit software package on bioactive molecules, we  
151 found that the ester functional group is the critical substructure that causes bioactivity on Cyp5  
152 (Figures 2a, b, and c) and ATXN (Figure 2d). Table S4 summarizes the substructures predicted  
153 to play a vital role in bioactivity toward nine different targets. The other 17 targets did not  
154 demonstrate as distinct clustering as the nine targets in Table S4 due to a relatively weak  
155 correlation between molecular structure and bioactivity.



156  
 157 **Figure 2:** Distribution of molecules in the CF dataset using semi-supervised metric learning. Each point  
 158 represents a molecule that is either bioactive (red circular edges) or inactive (light blue circular edges)  
 159 towards (a) CYP2C9, (b) CYP3A4, (c) CYP2D6, and (d) ATXN. The olive green-filled circles represent  
 160 molecules having the substructure depicted in the plot; i.e., (a, b) ester groups, (c) phenylpyrimidyl  
 161 groups, and (d) 4-benzyl-2-(4-fluorophenyl)-1,2-thiazole. The pink-filled circles in (c) represent  
 162 molecules with phenylethanone. The percentage value represents the ratio of the number of bioactive  
 163 molecules within the identified substructure. Table S3 lists the predicted substructures for specific  
 164 targets.  
 165



166  
167 We used structural alerts to cross-check the validity of the predicted substructures that play  
168 a crucial role in bioactivity. Within the bioinformatics community, structural alerts are  
169 molecular functional groups associated with a particularly adverse outcome, in our case,  
170 bioactivity.<sup>20,21</sup> We cross-referenced the ChEMBL dataset to our machine learning results since  
171 it contains structural alert information for some PFAS molecules.<sup>22</sup> Figure S3 shows structural  
172 alerts of the molecules that are bioactive on CYP2CP, and, as mentioned previously, the ester  
173 group was found to be the critical structure that induces interaction with Cyps.<sup>23,24</sup>

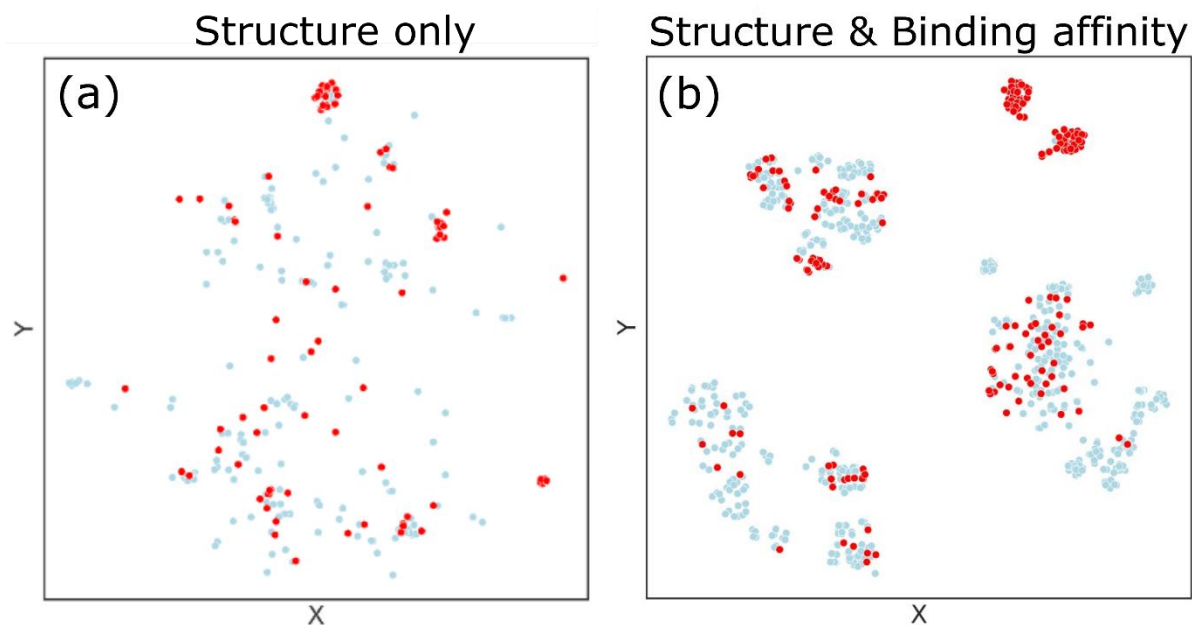
### 174 175 **3-3. Interactions between PFASs and targets**

176 We carried out molecular docking calculations with Autodock<sup>21</sup> to rationalize the  
177 underlying molecular causes of bioactivities in PFAS and predict their interaction with target  
178 enzymes. The Supporting Information gives additional details of our molecular docking  
179 calculations. We successfully docked all PFASs into the active sites of the targets and binned  
180 the binding affinity results based on their bioactivity with the target. Figure S5 displays one of  
181 the bioactive structures with the ester group of the CYP2C9-PFAS complex, methyl 4-[2-  
182 propyl-1-([4-trifluoromethyl]phenyl)sulfonyl]amino)-2-hexen-1-yl]benzoate.

183 To verify the correlation between the Autodock binding affinities and their bioactivity, we  
184 performed a dimension reduction procedure using unsupervised learning on the CF dataset,  
185 which consists of molecular structures with binding affinity data (see Figure 3). We used  
186 unsupervised learning here to make the point that unsupervised learning underperforms when  
187 only structural data is provided. Specifically, if the classification accuracy is improved with  
188 additional feature inputs, those features must contain some information to discriminate among  
189 the population.<sup>25,26</sup> In other words, if the inclusion of binding affinity data enhances the  
190 clustering accuracy, it provides another co-descriptor for bioactivity. Indeed, Figures 3b and  
191 3a show that descriptors consisting of chemical structures *and* binding affinity data give a better

192 separation/distinction between active and inactive molecules compared to the unsupervised  
193 learning results based only on chemical structures.

194



195 **Figure 3:** Clustering of molecules predicted with unsupervised learning (dimension reduction) on CF  
196 datasets containing (a) chemical structures and (b) chemical structures and binding affinities with  
197 CYP2C9. Each point represents a molecule that is either bioactive (red) or inactive (blue) towards  
198 CYP2C9.  
199

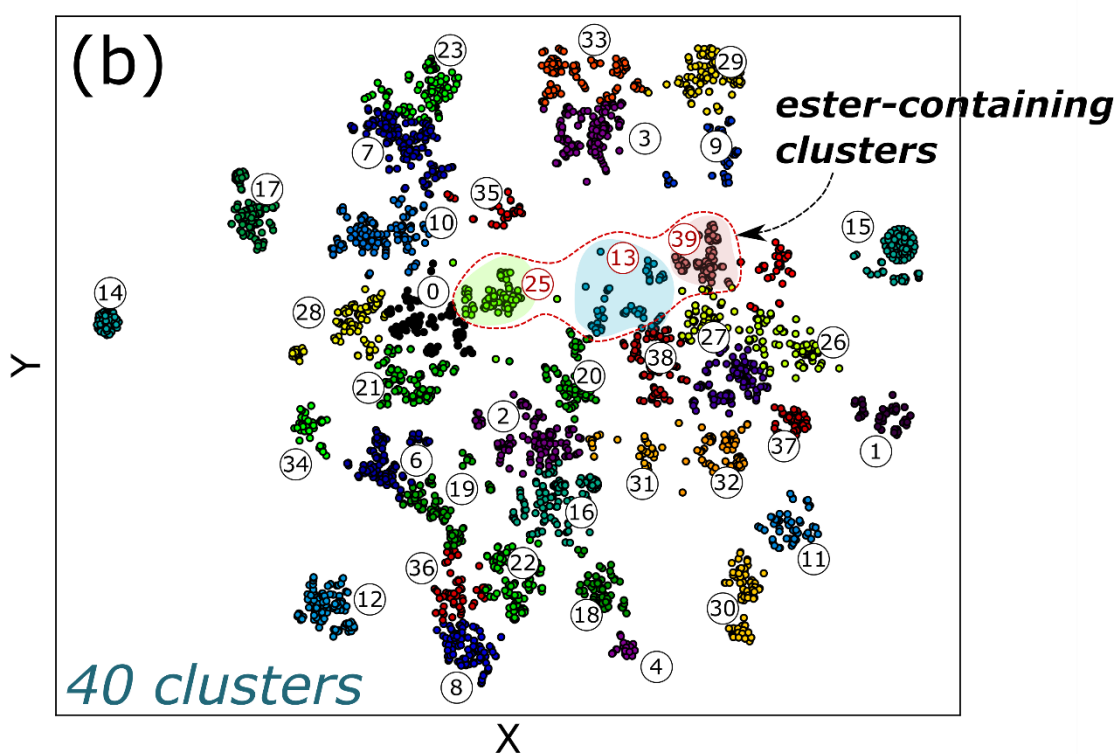
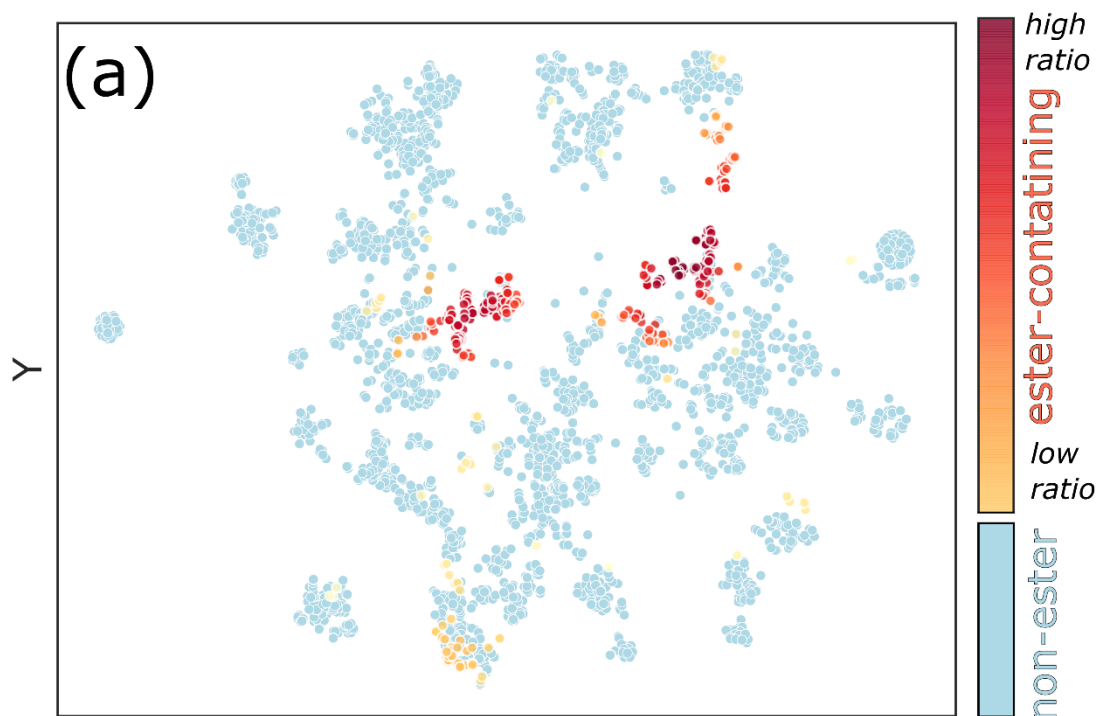
200

### 201 3-4. Bioactivity predictions on OECD dataset

202 In 2018, the Global Perfluorinated Chemicals Group<sup>27</sup> within the OECD published a list of  
203 4,730 PFASs to develop regulatory approaches for reducing the use of perfluorinated  
204 substances in products. However, researchers have yet to discover the bioactivities of the  
205 molecules in the list. Using the QSAR model developed in this work, we give predictions and  
206 a rationale for the bioactivities of molecules in the OECD list.

207 We performed molecular docking calculations on molecules containing the ester group  
208 among the OECD list to verify similar binding conformations. Of the 4,730 PFASs in the  
209 OECD list, 414 have an ester functional group. Figure S6 shows four different representative  
210 ester-containing molecules bound to CYP2C9. In particular, the ester-containing molecules in

211 the OECD list bind strongly with  $\text{Fe}^{2+}$  of the HEME group (an active site of Cyp enzyme),  
212 which is similar to the binding interactions that we observed in the CF dataset. Therefore, we  
213 expect a large portion of the 414 ester-containing molecules among the OECD list to form  
214 strong bonds with  $\text{Fe}^{2+}$  of the HEME group with a similar conformation, leading to bioactivity  
215 toward Cyp enzymes. Furthermore, based on our docking calculations, 87.7% of these 414  
216 molecules have a stronger binding affinity than -5 kcal/mol (the average binding affinity is -  
217 5.77 kcal/mol), which falls in the range of the mean binding affinity of the bioactive molecules  
218 from the CF dataset.



219  
 220 **Figure 4:** (a) OECD dataset classified by PC t-SNE and clustered based on the k-means clustering  
 221 method. The orange and yellow dots represent ester-containing molecules. The colors closer to red  
 222 (yellow) represent a higher (lower) concentration of bioactive molecules. (b) PFAS molecules included

223 in the OECD list are grouped into 40 clusters. Each point represents a molecule, and clusters 13, 25,  
224 and 39 denote a high ratio of ester-containing groups.

225  
226 We then clustered the OECD dataset into 40 clusters using the k-means clustering method.

227 Using both the clustered results (Figure 4b) and the distribution of ester-group-containing  
228 molecules (Figure 4a), we found that clusters 13, 25, and 39 contain ester functional groups.

229 Analyzing the CF dataset, we found that the ester group plays a possible role in bioactivity  
230 toward Cyp enzymes; that is, molecules in these clusters have a high probability of being  
231 bioactive against CYP2C9 and CYP3A4.

232 In summary, we have developed a new QSAR model validated with ChemLB structural  
233 alerts and molecular docking calculations, which constitutes the first application of semi-  
234 supervised metric learning for predicting/rationalizing bioactivities in PFASs. Using a semi-  
235 supervised metric learning algorithm, our machine-learning-based QSAR model accurately  
236 identified specific substructures, such as ester-containing groups, that play a possible role in  
237 determining bioactivities. With our semi-supervised learning approach, we obtained a distinct  
238 classification between bioactive and inactive molecules, resulting in an accuracy of up to 97.3%  
239 in the CF dataset. We also used semi-supervised metric learning to automatically  
240 classify/cluster and predict functional groups that could possibly play a role in bioactivity.

241 In addition, our machine learning model proposed a few significant substructures that could  
242 induce bioactivity, which were subsequently examined with molecular docking calculations.  
243 Most importantly, our machine learning predictions on bioactivities can provide a more  
244 efficient screening of potentially bioactive PFASs that can be used to complement *in vitro*  
245 assessments. All of our machine learning algorithms are publicly available (see Supporting  
246 Information), and we anticipate that researchers can further extend our methodology to screen  
247 other contaminants or analyze the potential bioactivity of PFAS molecules.

248

249 **Acknowledgments**

250 This material is based upon work supported by the National Science Foundation under grant  
251 No. CHE-1808242.

252

253 **Supporting Information**

254 Additional details on unsupervised and semi-supervised metric machine learning methods,  
255 additional details on molecular docking calculations, unsupervised machine learning results,  
256 and open-source Python codes for all the machine learning algorithms used in this work:  
257 [https://github.com/kha8128/PFAS\\_ML.git](https://github.com/kha8128/PFAS_ML.git). This information is available free of charge on the  
258 ACS Publications website.

259

260 **References**

261

- 262 (1) Hepburn, E.; Madden, C.; Szabo, D.; Coggan, T. L.; Clarke, B.; Currell, M.  
263 Contamination of Groundwater with Per- and Polyfluoroalkyl Substances (PFAS) from  
264 Legacy Landfills in an Urban Re-Development Precinct. *Environ. Pollut.* **2019**, *248*,  
265 101–113.
- 266 (2) Blake, B. E.; Pinney, S. M.; Hines, E. P.; Fenton, S. E.; Ferguson, K. K. Associations  
267 between Longitudinal Serum Perfluoroalkyl Substance (PFAS) Levels and Measures  
268 of Thyroid Hormone, Kidney Function, and Body Mass Index in the Fernald  
269 Community Cohort. *Environ. Pollut.* **2018**, *242*, 894–904.
- 270 (3) Guillette, T. C.; McCord, J.; Guillette, M.; Polera, M. E.; Rachels, K. T.; Morgeson,  
271 C.; Kotlarz, N.; Knappe, D. R. U.; Reading, B. J.; Strynar, M.; Belcher, S. M. Elevated  
272 Levels of Per- and Polyfluoroalkyl Substances in Cape Fear River Striped Bass  
273 (Morone Saxatilis) Are Associated with Biomarkers of Altered Immune and Liver  
274 Function. *Environ. Int.* **2020**, *136*, 105358.
- 275 (4) OECD. 033-066-C609-51.Pdf. *Series on Risk Management* **2018**, No. 39. (39), 1–24.
- 276 (5) Cousins, I. T.; Dewitt, J. C.; Glüge, J.; Goldenman, G.; Herzke, D.; Lohmann, R.;  
277 Miller, M.; Ng, C. A.; Scheringer, M.; Vierke, L.; Wang, Z. Strategies for Grouping  
278 Per- and Polyfluoroalkyl Substances (PFAS) to Protect Human and Environmental  
279 Health. *Environ. Sci.: Process. Impacts.* **2020**, *22*, 1444–1460.
- 280 (6) Hansch, Corwin.; Fujita, Toshio. P- $\sigma$ - $\pi$  Analysis. A Method for the Correlation of  
281 Biological Activity and Chemical Structure. *J. Am. Chem. Soc.* **2002**, *86*, 1616–1626.
- 282 (7) Cherkasov, A.; N. Muratov, E.; Fourches, D.; Varnek, A.; I. Baskin, I.; Cronin, M.;  
283 Dearden, J.; Gramatica, P.; C. Martin, Y.; Todeschini, R.; Consonni, V.; E. Kuz'min,  
284 V.; Cramer, R.; Benigni, R.; Yang, C.; Rathman, J.; Terflath, L.; Gasteiger, J.;  
285 Richard, A.; Tropsha, A. QSAR Modeling: Where Have You Been? Where Are You  
286 Going To? *J. Med. Chem.* **2014**, *57*, 4977–5010.

- 287 (8) Neves, B. J.; Braga, R. C.; Melo-Filho, C. C.; Moreira-Filho, J. T.; Muratov, E. N.;  
288 Andrade, C. H. QSAR-Based Virtual Screening: Advances and Applications in Drug  
289 Discovery. *Front. Pharmacol.* **2018**, *9*, 1275.
- 290 (9) Raza, A.; Bardhan, S.; Xu, L.; Yamijala, S. S. R. K. C.; Lian, C.; Kwon, H.; Wong, B.  
291 M. A Machine Learning Approach for Predicting Defluorination of Per- And  
292 Polyfluoroalkyl Substances (PFAS) for Their Efficient Treatment and Removal.  
293 *Environ. Sci. Technol. Lett.* **2019**, *6*, 624-629.
- 294 (10) Cheng, W.; Ng, C. A. Using Machine Learning to Classify Bioactivity for 3486 Per-  
295 and Polyfluoroalkyl Substances (PFASs) from the OECD List. *Environ. Sci. Technol.*  
296 **2019**, *53*, 13970–13980.
- 297 (11) Wang, Y.; Suzek, T.; Zhang, J.; Wang, J.; He, S.; Cheng, T.; Shoemaker, B. A.;  
298 Gindulyte, A.; Bryant, S. H. PubChem BioAssay: 2014 Update. *Nucleic Acids Res.*  
299 **2014**, *42*, 1075–1082.
- 300 (12) Rohrer, S. G.; Baumann, K. Maximum Unbiased Validation (MUV) Data Sets for  
301 Virtual Screening Based on PubChem Bioactivity Data. *J. Chem. Inf. Model.* **2009**,  
302 *49*, 169–184.
- 303 (13) Krewski, D.; Acosta, D.; Andersen, M.; Anderson, H.; Bailar, J. C.; Boekelheide, K.;  
304 Brent, R.; Charnley, G.; Cheung, V. G.; Green, S.; Kelsey, K. T.; Kerkvliet, N. I.; Li,  
305 A. A.; McCray, L.; Meyer, O.; Patterson, R. D.; Pennie, W.; Scala, R. A.; Solomon, G.  
306 M.; Stephens, M.; Yager, J.; Zeise, L. Toxicity Testing in the 21st Century: A Vision  
307 and a Strategy. *J. Toxicol. Environ. Health. B. Crit. Rev.* **2010**, *13*, 51–138.
- 308 (14) Subramanian, G.; Ramsundar, B.; Pande, V.; Denny, R. A. Computational Modeling of  
309  $\beta$ -Secretase 1 (BACE-1) Inhibitors Using Ligand Based Approaches. *J. Chem. Inf.*  
310 *Model.* **2016**, *56*, 1936–1949.
- 311 (15) Martins, I. F.; Teixeira, A. L.; Pinheiro, L.; Falcao, A. O. A Bayesian Approach to in  
312 Silico Blood-Brain Barrier Penetration Modeling. *J. Chem. Inf. Model.* **2012**, *52*,  
313 1686–1697.
- 314 (16) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model.* **2010**,  
315 *50*, 742–754.
- 316 (17) *RDKit*. <http://www.rdkit.org/> (accessed 2021-06-29).
- 317 (18) Rousseeuw, P. J. Silhouettes: A Graphical Aid to the Interpretation and Validation of  
318 Cluster Analysis. *J. Comput. Appl. Math.* **1987**, *20*, 53–65.
- 319 (19) Morris, G. M.; Ruth, H.; Lindstrom, W.; Sanner, M. F.; Belew, R. K.; Goodsell, D. S.;  
320 Olson, A. J. Software News and Updates AutoDock4 and AutoDockTools4:  
321 Automated Docking with Selective Receptor Flexibility. *J. Comput. Chem.* **2009**, *30*,  
322 2785–2791.
- 323 (20) Raies, A. B.; Bajic, V. B. In Silico Toxicology: Computational Methods for the  
324 Prediction of Chemical Toxicity. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* **2016**, *6*,  
325 147.
- 326 (21) Yang, H.; Lou, C.; Li, W.; Liu, G.; Tang, Y. Computational Approaches to Identify  
327 Structural Alerts and Their Applications in Environmental Toxicology and Drug  
328 Discovery. *Chem. Res. Toxicol.* **2020**, *33*, 1312–1322.
- 329 (22) Davies, M.; Nowotka, M.; Papadatos, G.; Dedman, N.; Gaulton, A.; Atkinson, F.;  
330 Bellis, L.; Overington, J. P. ChEMBL Web Services: Streamlining Access to Drug  
331 Discovery Data and Utilities. *Nucleic Acids Res.* **2015**, *43*, W612–W620.
- 332 (23) Cheng, X.; Klaassen, C. D. Perfluorocarboxylic Acids Induce Cytochrome P450  
333 Enzymes in Mouse Liver through Activation of PPAR- $\alpha$  and CAR Transcription  
334 Factors. *Toxicol. Sci.* **2008**, *106*, 29–36.
- 335 (24) Miners, J. O.; Birkett, D. J. Cytochrome P4502C9: An Enzyme of Major Importance in  
336 Human Drug Metabolism. *Br. J. Clin. Pharmacol.* **1998**, *45*, 525–538.

- 337 (25) Ashburner, J.; Klöppel, S. Multivariate Models of Inter-Subject Anatomical  
338 Variability. *Neuroimage* **2011**, *56*, 422–439.
- 339 (26) Chu, C.; Hsu, A. L.; Chou, K. H.; Bandettini, P.; Lin, C. P. Does Feature Selection  
340 Improve Classification Accuracy? Impact of Sample Size and Feature Selection on  
341 Classification Using Anatomical Magnetic Resonance Images. *Neuroimage* **2012**, *60*,  
342 59–70.
- 343 (27) *OECD Portal on Per and Poly Fluorinated Chemicals - OECD Portal on Per and Poly*  
344 *Fluorinated Chemicals*. [https://www.oecd.org/chemicalsafety/portal-perfluorinated-](https://www.oecd.org/chemicalsafety/portal-perfluorinated-chemicals/)  
345 [chemicals/](https://www.oecd.org/chemicalsafety/portal-perfluorinated-chemicals/) (accessed 2021-07-01).  
346