

# Exploring the Chemical Subspace of RPLC: a Data Driven Approach

Denice van Herwerden,<sup>\*,†</sup> Alexandros Nikolopoulos,<sup>†</sup> Leon P. Barron,<sup>‡,†</sup> Jake W. O'Brien,<sup>¶,†</sup> Bob W. J. Pirok,<sup>†</sup> Kevin V. Thomas,<sup>¶</sup> and Saer Samanipour<sup>\*,†,§,¶</sup>

<sup>†</sup>*Van 't Hoff Institute for Molecular Sciences (HIMS), University of Amsterdam, Amsterdam*

<sup>‡</sup>*MRC Centre for Environment and Health, Environmental Research Group, School of Public Health, Faculty of Medicine, Imperial College London, W12 0BZ, United Kingdom*

<sup>¶</sup>*Queensland Alliance for Environmental Health Sciences (QAEHS), The University of Queensland, Australia*

<sup>§</sup>*UvA Data Science Center, University of Amsterdam, Amsterdam*

E-mail: d.vanherwerden@uva.nl; s.samanipour@uva.nl

## Abstract

The chemical space is comprised of a vast number of possible structures, of which an unknown portion comprises the human and environmental exposome. Such samples are frequently analyzed using non-targeted analysis via liquid chromatography (LC) coupled to high-resolution mass spectrometry often employing a reversed phase (RP) column. However, prior to analysis, the contents of these samples are unknown and could be comprised of thousands of known and unknown chemical constituents. Moreover, it is unknown which part of the chemical space is sufficiently retained and eluted using RPLC. Therefore, we present a generic framework that uses a data driven approach to predict whether molecules fall ‘inside’, ‘maybe’ inside, or ‘outside’ of the

11 RPLC subspace. Firstly, three retention index random forest (RF) regression models  
12 were constructed that showed that molecular fingerprints are able to predict RPLC  
13 retention behavior. Secondly, these models were used to setup the dataset for building  
14 a RPLC RF classification model. The RPLC classification model was able to cor-  
15 rectly predict whether a chemical belonged to the RPLC subspace with an accuracy  
16 of 92% for the testing set. Finally, applying this model to the 91737 small molecules  
17 (i.e.,  $\leq 1000$  Da) in NORMAN SusDat showed that 19.1% fall outside of the RPLC  
18 subspace. Knowing which chemicals are outside of the RPLC subspace can assist in  
19 reducing potential candidates for library searching and avoid screening for chemicals  
20 that will not be present in RPLC data.

## 21 Introduction

22 The chemical space refers to a collection of all possible organic structures - for example,  
23 the GBD-17 database includes 116 billion possible organic molecules with a maximum of 17  
24 atoms, which is only a fraction of the chemical space.<sup>1-8</sup> Increasing the number of atoms  
25 only drastically increases these numbers and shows how vast the chemical space actually is.  
26 Even though these are possible structures, not all of them are likely to be present in the  
27 human and environmental exposome.<sup>8</sup> When evaluating the exposome, the main difficulty is  
28 that the contents of the samples taken are unknown prior to analysis and may comprise of  
29 thousands of both known and unknown constituents, particularly for small molecules (i.e.,  
30 molecular weight  $\leq 1000$  Da).<sup>9-16</sup> A frequently used approach for analyzing such samples  
31 is non-targeted analysis (NTA) via liquid chromatography (LC) coupled to high-resolution  
32 mass spectrometry (HRMS), for which a reversed phase (RP) LC selectivity is often used.<sup>8</sup>  
33 However, it is not yet known what part of the chemical space is covered by RPLC. The  
34 knowledge of the covered subspace also contains crucial information on chemicals that might  
35 not be visible in the final data even though they were present in the sample.<sup>3</sup>

37 Knowing what is separable with RPLC can have an improved outcome for both NTA  
38 and suspect screening. For NTA, the aim is to identify as much as possible of the potentially  
39 thousands of chemicals present in samples coming from, for example, biological or environ-  
40 mental backgrounds. Eliminating the potential candidates that fall outside of the chemical  
41 subspace of the selectivity (e.g., RPLC), reduces the number of false positive identifications.  
42 On the other hand, suspect screening is also a frequently used approach, where samples are  
43 screened for lists or even databases of compounds. Defining the subspace of a selectivity can  
44 reduce the number of potential candidates in these compound lists, reducing the computa-  
45 tional time required and the false positive matches with chemicals that cannot possibly be  
46 measured with this technique.

47

48 Separation data is usually limited to the mere assessment of whether the analyte retention  
49 time could fit in the range of the candidate's chemical class.<sup>17-20</sup> To take better advantage  
50 of the LC data, retention times are required to be initially converted to retention indices  
51 ( $r_i$ ), since the former are significantly influenced by the chromatography conditions, such as  
52 temperature, mobile phase composition, and gradients.<sup>20,21</sup> On the other hand,  $r_i$  values pro-  
53 vide a robust and highly reproducible way to express retention in liquid chromatography.<sup>20</sup>  
54 High reproducibility makes inter-laboratory results comparable, enabling both  $m/z$  and  $r_i$   
55 comparison with a reference and resulting in more confident suspect shortlisting.

56

57 As for any  $r_i$  system, different chromatography conditions should have negligible influence  
58 on the  $r_i$  value of the analytes, suggesting that there is a correlation between the  $r_i$  values  
59 and structural properties, expressed as molecular descriptors. This is the main principle  
60 used by the quantitative structure-retention relationship (QSRR) based models,<sup>22</sup> enabling  
61 the construction of QSRR models that either use all or a selection of descriptors to predict  $r_i$   
62 values.<sup>23-26</sup> However, difficulties arise when calculating descriptors due to convergence issues  
63 related to calculation time-out or local minima.<sup>25-27</sup> Moreover, descriptors can often be diffi-

64 cult to interpret, since they contain mathematical representations of the molecular structure.  
65 Alternatively, molecular fingerprints directly encode the molecular structure, making them  
66 more descriptive/understandable to interpret in relation to the chemical and do not require  
67 structural optimization (i.e., only uses 2D structural information), making them a potential  
68 alternative to descriptors.

69

70 In this paper, we present a data driven approach for a generic framework that enables  
71 quick screening of the RPLC chemical space, assuming that the molecules are in solution and  
72 can be injected into a system. A set of regression and classification models were built to assess  
73 whether a structure can theoretically be analyzed via RPLC. To build the RPLC classifica-  
74 tion model, firstly, we show the potential of using fingerprints for the prediction of  $r_i$  values  
75 for three retention index series, confirming that molecular fingerprints contain information on  
76 RPLC retention behavior. Three commonly used scales, namely: the n-alkylamide system,  
77 containing the n-alkylamide homologous series from n-propanamide to n-tetradecanamide  
78 (C3-C14)<sup>28</sup>, the  $r_i$  system developed by Aalizadeh et al. from the University of Athens re-  
79 ferred to as UoA, comprising of 18 reference compounds that were computationally selected in  
80 order to achieve a broad and reliable  $r_i$  reference system<sup>29</sup>, and the cocamide diethanolamine  
81 homologous series that is comprised of C(n = 0-23)-DEA chemicals<sup>30</sup> were employed for our  
82 model building. Secondly, we show the performance of the RPLC classification model and  
83 apply the model on a set of 91737 small molecules (i.e., molecular weight  $\leq$  1000 Da) from  
84 the NORMAN substance database (SusDat).

## 85 **Experimental Section**

### 86 **Overall Workflow**

87 The overall workflow for this work can be found in figure 1 and the details are explained  
88 in the following sections. In brief, a total of four random forest (RF) models were built, of

89 which three were  $r_i$  RF regression models (Figure 1A) and the fourth a RPLC RF classifi-  
90 cation model (Figure 1B). For building these models, a type of molecular fingerprint needed  
91 to be selected and the dataset obtained before model optimization and performance testing  
92 (Figure 1C). These models were used for evaluating the potential of using molecular finger-  
93 prints for prediction of retention behavior in RPLC and for setting up two of the classes  
94 for the fourth RF classification model. The latter refers to the ‘inside’ and ‘maybe’ inside  
95 class. Here, the ‘maybe’ class represents the chemicals that are poorly retained (i.e., close  
96 to  $t_0$ ) or require relatively high amounts of organic modifier to elute. All chemicals in be-  
97 tween the ‘maybe’ regions are classified as ‘inside’. For the RPLC classification model, a  
98 dataset with chemicals that were ‘inside’, ‘maybe’ inside, and ‘outside’ of the RPLC sub-  
99 space was constructed (Figure 1B). Finally, the application of the RPLC classification model  
100 was showcased by applying it on the NORMAN SusDat database, which is a collection of  
101 expert curated environmentally relevant chemicals that have been actively used for screening  
102 of complex samples. All datasets for constructing the models and the NORMAN SusDat  
103 database can be found on Figshare.<sup>31</sup>

104

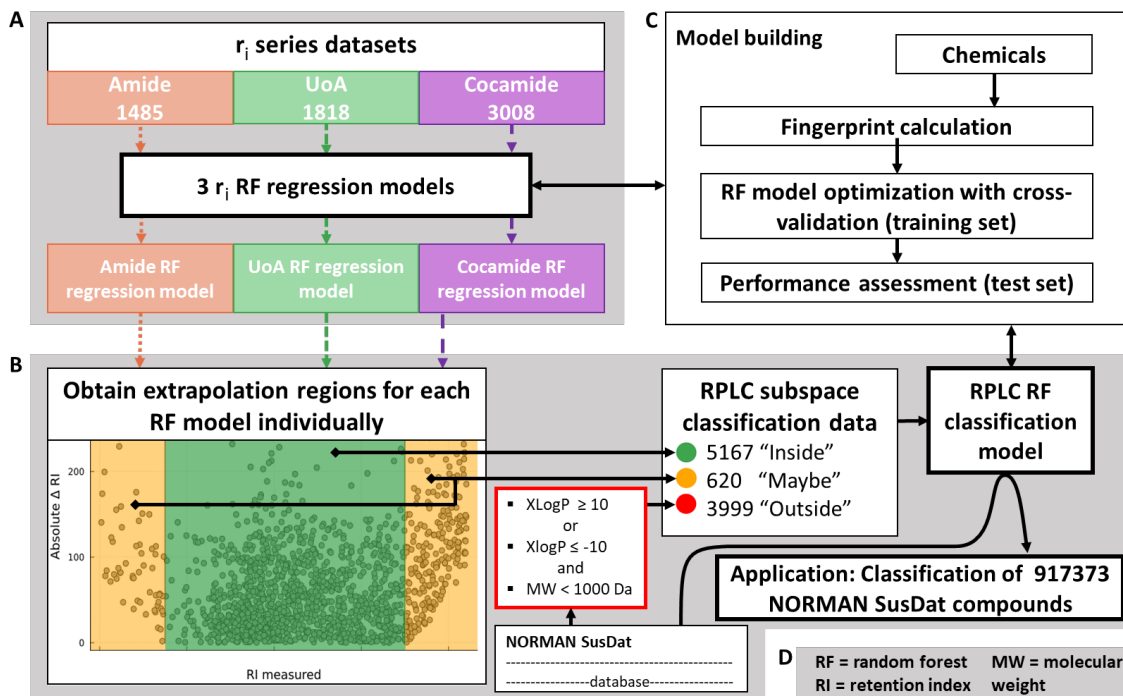


Figure 1: Workflow for construction of the RPLC classification model, comprising of the construction of three  $r_i$  RF regression models (A) and the construction the RPLC dataset for the RPLC RF classification model, which was applied to NORMAN SusDat (B). Finally C shows the model setup and D contains an overview of the abbreviations.

## 105 Fingerprint Calculations

106 The RF models were built using a combination of two different fingerprint series as inputs,  
 107 which included the AtomPairs2DFingerprintCount (2DAPC) and PubChem fingerprints,<sup>32</sup>  
 108 obtained through PaDEL.<sup>33</sup> The 2DAPC fingerprints counted the number of times two atoms  
 109 were present with a certain distance between themselves. For example, the molecule with  
 110 the SMILES 'NC(CC)CN' contains two times a distance of 3 between a C and N atom (i.e.,  
 111 C-x-x-N in the 2D molecular structure). The distances included ranges from 1 to 10 and  
 112 the elements considered were C, N, O, Cl, I, Br, F, P, S, Si, B, and X, where X represents  
 113 all halogens, yielding a total of 780 2DAPC fingerprints. As for the PubChem fingerprints,  
 114 only the portion of fingerprints containing ring information was used (i.e., PubChem finger-  
 115 print 115 - 262). These fingerprints were converted and reduced to a total of 10 additional

116 variables, which were the number of rings with a size of 3, 4, 5, 6, 7, 8, 9, 10, the number of  
117 aromatic rings, and the number of hetero-aromatic rings. Since the PubChem fingerprints  
118 are binary, there were multiple columns describing the same information but only differing  
119 in the number of a ring of a certain size. For example, for a ring size of 3, there were 2  
120 fingerprints, namely PubChem fingerprint 115 and 122, which were described as more than  
121 1 ring with a size of 3 or more than 2 rings with a size of 3, respectively. In case a molecule  
122 contained 2 rings with a size of 3, the PubChem fingerprints 115 would be 0 and 122 would  
123 be 1, which was converted to a single variable for our model containing the number of rings  
124 with a size of 3, meaning that this variable would be equal to 2 for this example case. An  
125 overview of which PubChem fingerprints were used for each of the 10 reduced PubChem  
126 variables can be found in table S2.

127

## 128 **Retention Index Random Forest Regression Models**

129 To show that fingerprints can be used to describe retention behavior in RPLC and for set-  
130 ting up the dataset for the RPLC classification model, random forest (RF) regression models  
131 were built using three different retention index series (Figure 1A). The three series used for  
132 this, were the amide<sup>28</sup>, University of Athens (UoA)<sup>29</sup>, and cocamide series.<sup>30</sup> For each of  
133 the series, the measured  $r_i$  were obtained from their respective articles, yielding 1485, 1818,  
134 and 3008 unique chemicals with measured  $r_i$  values for the amide, UoA, and cocamide series,  
135 respectively. For all chemicals, the 2DAPC and PubChem fingerprints were calculated ac-  
136 cording to Section ‘Fingerprint Calculations’. For each  $r_i$  series, data was split into a training  
137 and test set, at random, with a ratio of 0.85:0.15, ensuring similar coverage of the  $r_i$  range  
138 in both sets. The test set was only used for testing and thus never used for training. For  
139 optimization of the RF regression models, the training set was used with a 0.8:0.2 split for  
140 training and cross-validation, respectively. This ratio of split has been shown to be effective  
141 in such data sets.<sup>25,26,34,35</sup> The RF regression models used a third of the features (i.e., 264)

142 for training each tree. The parameters that were optimized were the minimum number of  
143 samples per leaf and the number of trees. The minimum number of samples per leaf tested  
144 were 4, 6, 8, 10, 15, and 20. The tested number of trees were 50, 100, 150, 200, 250, 300,  
145 350, 400, 500, 600, 700, 800, 900, and 1000. In addition, the random state for splitting the  
146 cross-validation set and selection of the features in the RF models for each tree was also  
147 varied with values of 1, 2, and 3. The accuracy of the cross-validation set for each possible  
148 combination of the minimum number of samples per leaf, number of trees, and random state  
149 was used for the optimization of the RF models. After obtaining the optimized models for  
150 the amide, UoA, and cocamide series, the applicability domains were assessed according to  
151 Section ‘Applicability Domain Calculations’. Finally, for each  $r_i$  series, the optimized model  
152 and applicability domain assessment were applied on the test set to evaluate the performance  
153 of the model on unseen data.

154

## 155 **RPLC Random Forest Classifier**

156 The dataset for building the RPLC classifier model was comprised of three classes: ‘inside’,  
157 ‘maybe’, and ‘outside’ the RPLC subspace (Figure 1B). The ‘outside’ chemicals were ob-  
158 tained from the NORMAN SusDat database based on their extreme XLogP values. Here,  
159 the XLogP was chosen rather than the logD due to the fact that it is easier to predict, more  
160 stable, and more accurate.<sup>36</sup> For the ‘outside’ case, a total of 3999 compounds with a XLogP  
161 value above 10 or below -10 and with a molecular weight below 1000 Da were obtained. As  
162 for the ‘inside’ and ‘maybe’ chemicals, these were obtained from the experimentally defined  
163  $r_i$  values by the three  $r_i$  series. For each of the series, the absolute difference between the  
164 predicted and measured  $r_i$  (i.e., the residuals) versus the measured  $r_i$  values were plotted  
165 and the regions of extrapolation were identified. These regions were obtained based on the  
166 increasing residuals that were caused by the inherent over estimation and under estimation  
167 of a RF regression model, which are associated with either extremely low or extremely high



168  $r_i$  values, respectively. These regions correspond to chemicals that elute close to  $t_0$  or are  
169 very difficult to elute from the column (i.e., require a relatively high percentage of organic  
170 modifier). The chemicals with a measured  $r_i$  in these extrapolation regions were labeled  
171 as ‘maybe’ and the remaining chemicals were labeled as ‘inside’ the RPLC subspace. This  
172 yielded a total of 620 ‘maybe’ and 5167 ‘inside’ compounds. Whenever a chemical was found  
173 in multiple classes (i.e., it was present in multiple datasets of the  $r_i$  models), it was removed  
174 from the lower ranking RPLC classes and kept in the highest ranking RPLC class (i.e., ‘in-  
175 side’ > ‘maybe’ > ‘outside’ RPLC class rank). For example, if a chemical was found in the  
176 ‘maybe’ region for UoA and in the ‘inside’ for Cocamide, it would be classified as ‘inside’.  
177 More details on the division between the ‘inside’ and ‘maybe’ classification can be found  
178 in Section ‘RPLC Classification Model’ as these are based on the results of the three RF  
179 regression models.

180

181 The dataset described above was used for building the RPLC classifier model with a  
182 training set/test set split of 0.85:0.15, ensuring equal distribution of each class in both sets.  
183 The optimized RF classifier model was obtained using the same approach as for the RF re-  
184 gression models (see Section ‘Retention Index Random Forest Regression Models’). For this  
185 model, the applicability domain was also obtained as described below. Finally, the optimized  
186 RPLC classification model and applicability domain assessment was applied to the test set  
187 and the performance was evaluated.

188

## 189 **RPLC Space Prediction for NORMAN SusDat**

190 To showcase the model’s potential, it was applied to the NORMAN SusDat database.<sup>5</sup> For  
191 this, the 2DAPC and reduced PubChem fingerprints for a total of 91737 chemicals with a  
192 molecular weight below 1000 Da from SusDat were calculated. These fingerprints were then  
193 used to calculate the leverage of each chemical with the RPLC classifier training set, as

194 explained in the next section ‘Applicability Domain Calculations’, and to apply the RPLC  
195 classifier model to each of the SusDat chemicals. To visualize the coverage of each class  
196 (i.e., ‘inside’, ‘maybe’, and ‘outside’ the RPLC subspace), the molecular weight was plotted  
197 against the XLogP, which were obtained from the descriptor calculations of PaDEL.

198

## 199 **Applicability Domain Calculations**

200 Applicability domain calculations were used to assess whether the training data, used in the  
201 random forest models, sufficiently covered the variable space for new chemicals on which the  
202 models need to be applied.<sup>25,37</sup> This was done through leverage calculations of a chemical  
203 with the entire training set, yielding a distance of that chemical to the training set. Equation  
204 1 shows how the leverage is calculated, where  $X$  is the training data matrix and  $x_i$  is the  
205 sample vector, both containing the 2DAPC and reduced PubChem fingerprints for our mod-  
206 els. To set a threshold for this, the leverage was calculated for all training samples with the  
207 entire training set of a model, yielding values between 0 and 1. Then, a leverage threshold  
208 was obtained that covered 95% of the training data. If a chemical, compared to the training  
209 set of the model in question, had a value lower than the leverage threshold, the compound  
210 was within the applicability domain, and, if the value was above the leverage threshold, the  
211 results should be taken with care as the training data might not be sufficiently describing  
212 the variable space for the new compound.

213

$$l_{ii} = x_i(X^T X)^{-1}x_i \quad (1)$$

## 214 **Calculations and Code Availability**

215 The calculations and development of the models were executed on a personal computer with  
216 12 CPUs and 32 GB of RAM, using Windows 10. The  $r_i$  regression and RPLC classifi-

217 cation models were developed and evaluated with the Julia programming language (v1.6).  
218 The code for using the  $r_i$  regression models and RPLC space prediction model is available  
219 at: [https://bitbucket.org/Denice\\_van\\_Herwerden/riprediction/src/main/](https://bitbucket.org/Denice_van_Herwerden/riprediction/src/main/). This Ju-  
220 lia package contains functions for obtaining the required 2DAPC and reduced PubChem  
221 fingerprints and for using the  $r_i$  regression models and RPLC sub space classification model.  
222

## 223 Results and discussion

### 224 Retention Index Random Forest Regression Models

225 All three  $r_i$  regression models obtained an accuracy of 81% for the training set and, for the  
226 test set. The amide, UoA, and cocamide models had an accuracy of 68%, 70%, and 67%,  
227 respectively. The  $r_i$  regression models were built and optimized for the amide, UoA and  
228 cocamide series. Grid optimization of each of these models showed that the number of trees  
229 did not influence the performance of the model (Figures S1, S2, and S3). Therefore, to  
230 keep the model light, 200 trees were selected. As for the minimum number of samples per  
231 leaf, 8 was found to be the optimum, based on the training and cross-validation accuracy.  
232 When evaluating the predicted versus the measured  $r_i$  values for these models a trend of over  
233 prediction for lower  $r_i$  values and under prediction of higher  $r_i$  values was found (Figures S4,  
234 S6, and S8), corresponding to the regions where the RF regression models were extrapolat-  
235 ing. These regions were used for establishing the ‘maybe’ areas for the RPLC classification  
236 dataset.

237

238 Most compounds (i.e., 88.5%) in our test set appeared to be within the applicability  
239 domain of each model. To obtain the applicability domains of these models, a 95% leverage  
240 threshold of 0.189 for amide, 0.652 for UoA, and 0.424 for cocamide was found for the train-  
241 ing sets. For the training set the leverage values range between 0 and 1, meaning that the

242 lower threshold for the amide model showed how similar most of the amide compounds were  
243 to each other, while for the UoA and cocamide models, the higher thresholds corresponded  
244 with the larger variety of chemical structures found in the dataset. When the leverage cal-  
245 culations were applied on the test sets for these models, a total of 22, 34, and 54 compounds  
246 were found to be outside of the applicability domain for the amide, UoA, and cocamide  $r_i$   
247 models, respectively. This does not necessarily mean that the predicted outcome for these  
248 cases was wrong, as can be seen in figures S4, S6, and S8. Here, most chemicals outside the  
249 applicability domain still follow the trend of the other data points. However, the outcome  
250 should be taken with care as the model might insufficiently cover the chemical space for a  
251 new compound in question, especially for leverage values  $> 1$ . It should be noted that the  
252 largest training set leverage value obtained from our applicability domain calculations was 1.

253

254 The cocamide RF regression model used the most fingerprints for the prediction of the  
255  $r_i$  indices (i.e., 215 fingerprints), while the UoA and amide  $r_i$  models used 165 and 61, re-  
256 spectively. The low number of fingerprints used for amide was not surprising due to the  
257 fact that the compounds in this  $r_i$  series are only comprised of C, H, N, and O. Hence, the  
258 amide  $r_i$  model only used the 2DAPC fingerprint counts with a certain distance between C,  
259 N, and O atoms. At first sight, this was also noticeable when comparing the top 20 most  
260 important fingerprints for the three  $r_i$  models (S3). The most contributing fingerprints for  
261 the amide  $r_i$  model were the distances 1 till 7 between two C atoms with importance ranging  
262 between 27% and 4%. As for the UoA  $r_i$  model, C-Cl and C-X distance begin to contribute  
263 more to the model and the most important fingerprint (i.e., distance 7 between C-C) only  
264 contributes 9.6%, having an overall more divided importance between a larger group of con-  
265 tributing features than the amide model. Finally, a similar trend was also observed for the  
266 cocamide model, except that the C-X distances start to play a more important role than the  
267 C-Cl distances, which could be explained by the higher number of halogens present in the  
268 compounds from the cocamide dataset. This variability in important features used in each

269  $r_i$  regression model shows that different structures may be better captured by one  $r_i$  model  
270 vs another, due to the diversity of training set in terms of chemical structures. This, also,  
271 further indicates the need for a more generic model incorporating the information from all  
272 three  $r_i$  models.

273

274 Overall, these models show that a combination of the 2DAPC fingerprints and the re-  
275 duced PubChem fingerprints can be used to predict  $r_i$  values. All three models performed  
276 almost equally well with negligible deviations for the training set accuracy. However, de-  
277 pending on the chemicals for which  $r_i$  would be predicted, it is advised to evaluate which  
278 model would be most suitable based on the leverage applicability domain calculations.

279

## 280 RPLC Classification Model

281 To build the RPLC classification model, it was assumed that the chemicals are in solution  
282 and that the chemicals can be injected into a system. Additionally, the model focuses on  
283 whether an analyte could be analyzed with RPLC regardless of experimental parameters or  
284 sample pretreatment. The dataset for this was comprised of 5167 ‘inside’, 620 ‘maybe’ in-  
285 side, and 3999 ‘outside’ chemicals for the RPLC subspaces. The ‘outside’ cases were obtained  
286 from NORMAN SusDat with extreme XLogP values, while the ‘inside’ and ‘maybe’ cases  
287 came from the three  $r_i$  regression models. In figures S10, S11, and S12 the extrapolation  
288 limits for each of the models are defined. For  $r_i$  range for the ‘inside’ RPLC subspace for  
289 the amide, UoA, and cocamide series were 350-900, 100-900, and 250-1300, respectively. All  
290 compounds that had a higher or lower  $r_i$  value for the corresponding range of the model it  
291 was coming from, were classified as ‘maybe’ inside the RPLC subspace, due to the fact that  
292 these chemicals either elute close to  $t_0$  or require high percentages of organic eluent to be  
293 eluted.

294

295 The final optimized classification model resulted in an accuracy of 94% and 92% for the  
296 training and test set, respectively (Figures 2, and S15). In this case 200 trees and 8 minimum  
297 samples per leaf was found to be the optimum for the model (Figure S13). For the training  
298 and test set, 90.8% and 87.7% of the ‘inside’ and ‘maybe’ cases were correctly classified, 7.4%  
299 and 9.3% of the ‘inside’ and ‘maybe’ cases were wrongly classified as a ‘maybe’ or ‘inside’  
300 case, respectively, and 1.7% and 3.0% of the ‘inside’ and ‘maybe’ cases were wrongly classi-  
301 fied as ‘outside’. For the ‘outside’ cases, 0.7% and 1.5% of the cases were wrongly classified  
302 as an ‘inside’ or ‘maybe’ case and 99.3% and 98.5% of the cases was correctly classified as  
303 an ‘outside’ case for the training and test set, respectively. Overall, considering that the  
304 wrongly classified ‘inside’ and ‘maybe’ cases as ‘maybe’ and ‘inside’, respectively, still are  
305 considered part of the RPLC subspace, the performance of the model was very good with  
306 only 2.4% of all cases being wrongly classified as ‘inside’ or ‘maybe’ while being an ‘outside’  
307 or vice versa for the test set.

308

309 As for the applicability domain of the RPLC classification model, the 95% leverage  
310 threshold of the training set was 0.209 (Figure S14). In total, 102 compounds from the test  
311 set (i.e., 6.9%) had a leverage with the training set that was higher than 0.209, of which 31  
312 cases had leverage values above 1. Out of these 102 cases only 10 were wrongly classified  
313 and had leverage values ranging between 0.209 to the most extreme (i.e., 809.255), showing  
314 that in this case higher leverage values did not necessarily mean that the model would have  
315 a higher error. However, it should be noted that cases with a very large leverage should be  
316 considered with extra care, as they may have a higher level of uncertainty.

317

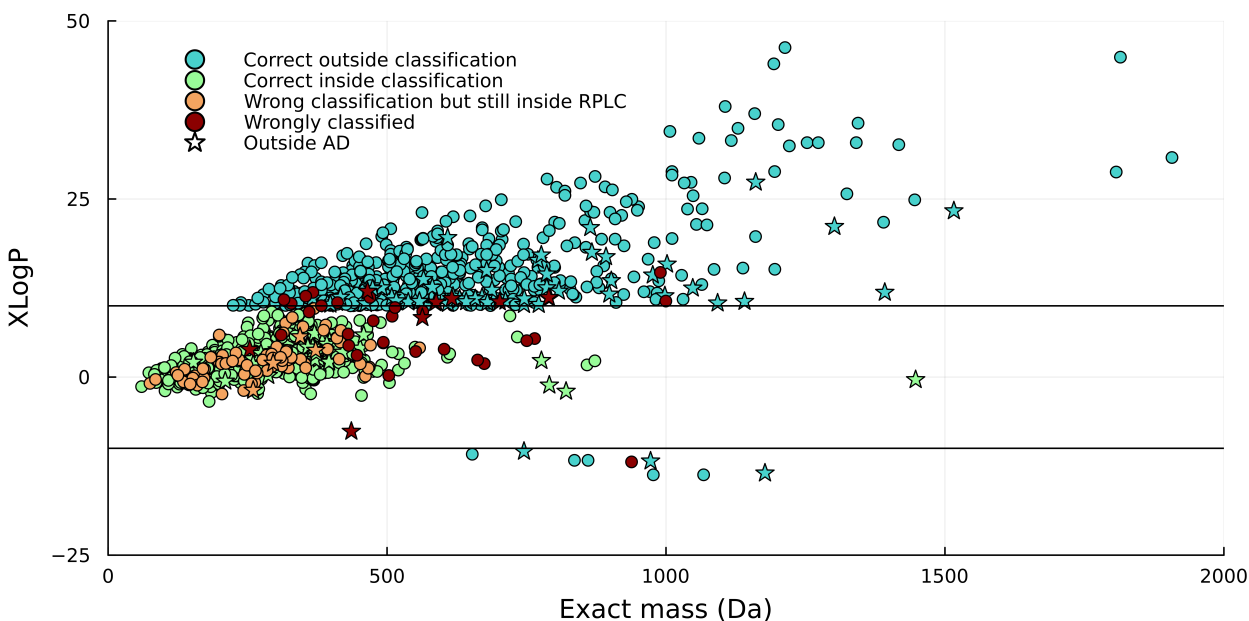


Figure 2: XLogP values versus the molecular weight for the RPLC classification test set. In blue are the correctly classified ‘outside’ cases, in green are the correctly classified ‘inside’ and ‘maybe’ cases, in orange are the wrongly classified ‘inside’ cases as ‘maybe’ and vice versa, in red the wrongly classified ‘inside’ and ‘maybe’ cases as ‘outside’ and the wrongly classified ‘outside’ cases as ‘inside’. The star markers show the compounds that were outside the 95% applicability domain of the RPLC classification training set

318 A total of 280 features were contributing to the RPLC classification model. This is more  
 319 than for each of the three  $r_i$  regression models, which was expected due to the higher variety  
 320 in chemical structures used in the RPLC classification model. The 20 most contributing fea-  
 321 tures are mainly described by ring related features and distances between combinations C,  
 322 N, and O atoms. A previous version of the model that was tested, using only the 2DAPC fin-  
 323 gerprints, frequently wrongly classified ‘inside’ as ‘outside’ due to the high degree of cyclicity  
 324 in the chemical structures (e.g., InChIKey: IUKLSMSEHKDIIP-BZMYINFQSA-N). Hence,  
 325 the addition of the reduced PubChem fingerprints better captures these chemical properties.  
 326 As a result, the number of rings with a size of 6, the minimum number of aromatic rings, and  
 327 the number of rings with a size of 5 were also part of the top 20 most contributing features.

328

329 In total, considering the extreme misclassifications, 9 out of 599 ‘outside’ chemicals were

330 wrongly classified as ‘inside’ or ‘maybe’ inside the RPLC subspace and 14 out of the 767  
331 ‘inside’ and 12 out of the 102 ‘maybe’ cases were classified as ‘outside’ the RPLC subspace.  
332 Two of the nine wrongly classified ‘outside’ cases were organic complexes that, in the mobile  
333 phase, would be analyzed as multiple smaller molecules (e.g., Gadopentetic acid dimeglu-  
334 mine salt). Also, another case was a surfactant containing a positive and negative charge  
335 (i.e., 4-Dodecyl-2-[(2-nitrophenyl)azo]phenol). This case was a chemical that falls ‘outside’  
336 of the RPLC space due to its predicted XLogP value of 10.452. However, the charges on  
337 this molecule would make it difficult to calculate this value accurately. Lexidronam was one  
338 of the ‘maybe’ cases that was classified as ‘outside’, due to a large leverage value of 26.0  
339 and the fact that it elutes at  $t_0$  (i.e., amide scale  $r_i$  of 206 versus urea  $r_i = 200$ ), indicating  
340 the need for special gradients to be able to retain such a chemical. As for the ‘inside’ cases  
341 that were wrongly classified as ‘outside’, generally larger, branched (e.g., SCHEMBL312614),  
342 or hydrolyzing (e.g., Bis[2-(perfluorohexyl)ethyl] Phosphate, respectively) chemicals showed  
343 higher likelihood of such misclassifications. Again these are structures that may require very  
344 specific adjustment of experimental condition (e.g., pH of mobile phase) to fit them within  
345 the RPLC analyzable chemical subspace.

346

347 Overall, our RPLC classification model was highly successful in identifying the chemical  
348 structures that are easily analyzable via RPLC (i.e., ‘inside’ cases) as well as the ‘maybe’ and  
349 ‘outside’ cases. The classification model used a combination of similar molecular fingerprints  
350 as those used by the three  $r_i$  models, taking advantage of all the structural information.

## 351 **NORMAN SusDat Chemical Space Prediction**

352 Finally, the RPLC classification model was applied to a set of small molecules (i.e., molecular  
353 weight < 1000) from the NORMAN SusDat database. In total, 80503 chemicals were within  
354 the applicability domain with leverage values  $\leq 0.209$ , 6570 compounds had leverage values  
355 between 0.209 and 1, and 4664 compounds had even larger leverages. This showed that the



356 RPLC classification model was suitable for a large variety, 87.8%, of compounds present in  
357 SusDat. The model predicted that 79.0% of the compounds would fit ‘inside’ the RPLC  
358 subspace, 2.0% was ‘maybe’ in this space, and 19.1% was ‘outside’ of the RPLC subspace.  
359 Examples of molecules classified as ‘inside’, ‘maybe’, and ‘outside’ were carbamazepine, su-  
360 dan I, and coronene, respectively. When comparing the relationship between XlogP and  
361  $r_i$ , it is clearly observable that these parameters, even though relatively linear, are insuffi-  
362 cient to determine if a chemical fits the RPLC subspace, figure 3. In figures S16,S17, and  
363 S18, the XlogP values of the chemicals with the same  $r_i$  range vary between -10 to +10 units.

364

365 Using the developed classification models implies that for screening RPLC samples against  
366 databases such as SusDat, 1/5 of the overall time can be saved, which becomes even more  
367 significant when applying it to larger sample sets. Additionally, this will result in higher  
368 confidence identifications when performing database matching for an RPLC NTA method  
369 with SusDat, by reducing the overall number of potential candidates and thus false positive  
370 identifications.

371

372 The amide  $r_i$  model is the least suited scale based on its applicability domain coverage  
373 since only 44500 (i.e., 48.5%) chemicals fell within the applicability domain. For the chem-  
374 icals that were outside the applicability domain, 18988 had a leverage value between 0.189  
375 and 1 (i.e., similar to the full training set) and 28249 had an even higher leverage value. As  
376 for the UoA and cocamide  $r_i$  models, 71022 (i.e., 77.4%) and 74252 (i.e., 80.9%) compounds  
377 were within the applicability domain. For the UoA model, 3421 and 17294 chemicals had a  
378 leverage value below and above 1, respectively, and the cocamide model had 5947 chemicals  
379 with a leverage value below 1 and 11538 chemicals with higher leverage values. Figures S16,  
380 S17, and S18 show the coverage of the ‘inside’, ‘maybe’, and ‘outside’ RPLC classes in terms  
381 of the XLogP values versus the predicted  $r_i$  values for the amide, UoA, and cocamide series.  
382 As expected the chemicals classified as ‘maybe’ inside RPLC are mainly clustering around

383 the lower and higher  $r_i$  values. While the chemicals classified as ‘outside’ the RPLC space  
384 span the entire  $r_i$  range for each of the three  $r_i$  series, suggesting that  $r_i$  prediction would  
385 also be insufficient to define the boundaries of the RPLC subspace.

386

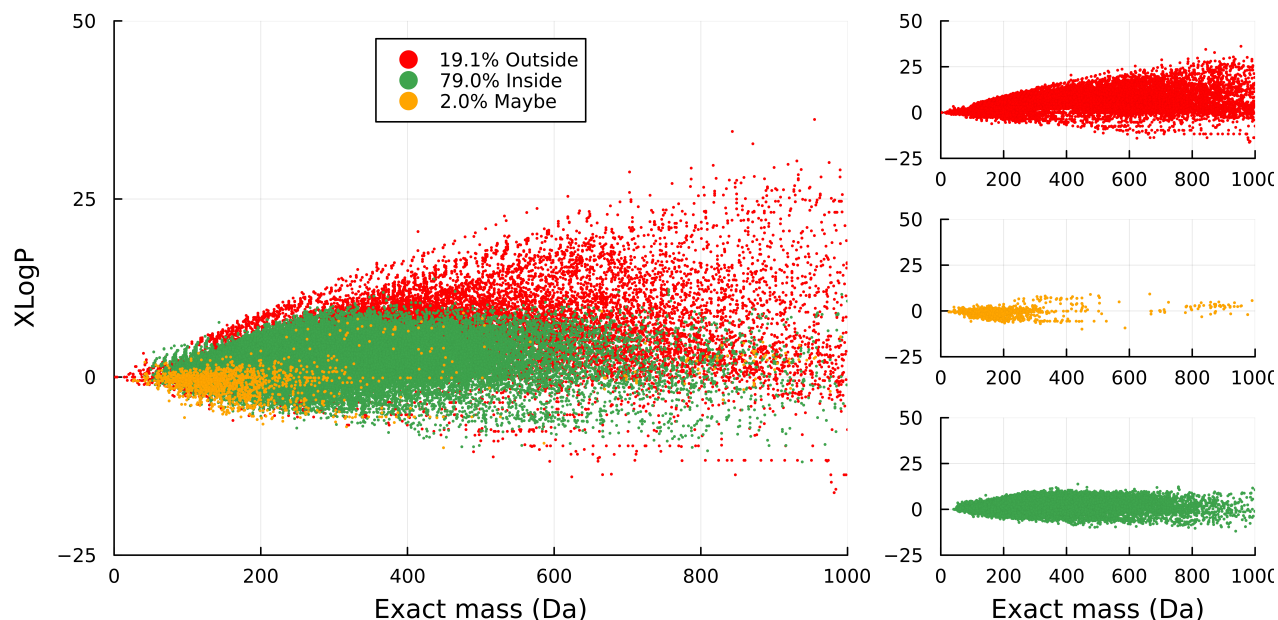


Figure 3: XLogP values versus the molecular weight for the NORMAN SusDat database compounds with a molecular weight below 1000 Da. In red, orange, and green are the compounds that were classified as ‘outside’, ‘maybe’, and ‘inside’ the RPLC chemical space, respectively. The subplots on the left show the coverage of the individual classes.

## 387 Potentials and Limitations

388 Overall, we developed four models for exploration of the RPLC subspace. The  $r_i$  regression  
389 models showed that fingerprints can be used for describing RPLC retention indices. Con-  
390 sequently, these fingerprints were used for RPLC classification model building. This model  
391 was able to predict whether chemicals were ‘inside’, ‘maybe’ inside, or ‘outside’ of RPLC  
392 chemical subspace with an accuracy of 92% on the test set. Applying the RPLC classification  
393 model on NORMAN SusDat showed that 19.1% of the compounds were classified as ‘out-  
394 side’ the RPLC subspace. This means that, when performing identification on NTA RPLC

395 samples, candidates classified as ‘outside’ compounds are unlikely to be the true structure of  
396 the chemical and can be removed to reduce the number of false positive identifications. In  
397 terms of suspect screening, it can save computational time since the ‘outside’ chemicals fall  
398 ‘outside’ of the RPLC subspace and thus should not be screened for. Additionally, 87.8%  
399 of NORMAN SusDat was within the applicability domain of the RPLC classifier, showing  
400 good coverage of a variety of compounds. The RPLC classification model also showed that  
401 the XLogP or  $r_i$  values alone are not sufficient to define the RPLC subspace.

402

403 The RPLC classification model overall did have more difficulties with regard to more  
404 bulky and branched or surfactant-like chemicals. Additionally, the model was not able to  
405 properly predict the RPLC subspace class of chemicals that are organic complexes, due to  
406 the fact that in solution those are dissociated into multiple individual structures. The latter  
407 is not a major limitation for the model itself, since, using expert knowledge, they can be  
408 easily identified. Generally, as knowledge on analyzable chemicals with RPLC grows, the  
409 model could easily be rebuilt and expanded for the range of analytes. Ideally, when sufficient  
410 data becomes available, selectivity classification models could be constructed for other se-  
411 lectivities (e.g., HILIC). This allows for further understanding of what part of the chemical  
412 space is actually covered by the selectivities used in NTA and what we are missing.

413

414 Moreover, the RPLC classification model uses a data driven approach and is intended  
415 for quick screening of the RPLC chemical space. The model assumes that compounds are  
416 analyzable with RPLC regardless of the chemicals solubility, experimental parameters, or  
417 pretreatment steps taken. This means that it cannot be assumed that chemicals ‘inside’ the  
418 RPLC space will be analyzable with every RPLC method. Here, the method subspace plays  
419 a major role when looking at what individual NTA methods can cover, becoming an even  
420 more complex issue due to the fact that sample pretreatment, gradient program’s, and RP  
421 column selectivities play a large influence on this. Defining the method chemical space would

422 be the next step in understanding what part of the vast chemical space we are covering and,  
423 more importantly, excluding with our current NTA methods.

424

## 425 **Acknowledgement**

426 The authors thank the Environmental Monitoring and Computational Mass Spectrometry  
427 (www.emcms.info) group for their insights and feedback. The Queensland Alliance for En-  
428 vironmental Health Sciences. Finally, the University of Queensland gratefully acknowledges  
429 the financial support from the Queensland Department of Health. J.W.O is the recipient of  
430 an NHMRC Emerging Leadership Fellowship (EL1 2009209).

## 431 **Supporting Information Available**

432 Overview of performance for using different types of molecular fingerprints, composition of  
433 reduced PubChem fingerprints, optimization, prediction, leverage, and feature importance  
434 results for the 3 RF regression models and the RPLC classification model, and the RPLC  
435 classification of NORMAN SusDat visualized by plotting the XLogP values versus the pre-  
436 dicted  $r_i$  values for the three  $r_i$  regression models.

## 437 **Author Information**

438 Corresponding Author:

439 Saer Samanipour

440 Van 't hoff institute for molecular sciences (HIMS),

441 University of Amsterdam,

442 the Netherlands

443 Email: s.samanipour@uva.nl

444

445 Denice van Herwerden

446 Van 't hof institute for molecular sciences (HIMS),

447 University of Amsterdam,

448 the Netherlands

449 Email: d.vanherwerden@uva.nl

450

## References

- (1) Ruddigkeit, L.; Deursen, R. V.; Blum, L. C.; Reymond, J. L. Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17. *Journal of Chemical Information and Modeling* **2012**, *52*, 2864–2875.
- (2) Reymond, J. L. The Chemical Space Project. *Accounts of Chemical Research* **2015**, *48*, 722–730.
- (3) Black, G. et al. Exploring chemical space in non-targeted analysis: a proposed ChemSpace tool. *Analytical and Bioanalytical Chemistry* **2023**, *415*, 35–44.
- (4) Williams, A. J.; Grulke, C. M.; Edwards, J.; McEachran, A. D.; Mansouri, K.; Baker, N. C.; Patlewicz, G.; Shah, I.; Wambaugh, J. F.; Judson, R. S.; Richard, A. M. The CompTox Chemistry Dashboard: a community data resource for environmental chemistry. *J. Cheminf.* **2017**, *9*, 61.
- (5) NORMAN SusDat. <https://www.norman-network.com/nds/susdat/>.
- (6) Lipinski, C. A.; Dominy, B. W.; Feeney, P. J. drug delivery reviews Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. 1997.
- (7) Bohacek, R. S.; McMartin, C.; Guida, W. C. The art and practice of structure-based drug design: A molecular modeling perspective. 1996.
- (8) Hulleman, T.; Turkina, V.; O, J. W.; Chojnacka, A.; Thomas, K. V.; Samanipour, S. Critical assessment of covered chemical space with LC-HRMS non-targeted analysis.
- (9) Schulze, B.; Jeon, Y.; Kaserzon, S.; Heffernan, A. L.; Dewapriya, P.; O'Brien, J.; Gomez Ramos, M. J.; Ghorbani Gorji, S.; Mueller, J. F.; Thomas, K. V.; Samanipour, S. An assessment of quality assurance/quality control efforts in high resolution mass spec-

- 474        trometry non-target workflows for analysis of environmental samples. *TrAC Trends in*  
475        *Analytical Chemistry* **2020**, *133*, 116063.
- 476 (10) Schymanski, E. L. et al. Non-target screening with high-resolution mass spectrometry:  
477        Critical review using a collaborative trial on water analysis. *Anal. Bioanal. Chem.* **2015**,  
478        *407*, 6237–6255.
- 479 (11) Werner, E.; Heilier, J.-F.; Ducruix, C.; Ezan, E.; Junot, C.; Tabet, J.-C. Mass spec-  
480        trometry for the identification of the discriminating signals from metabolomics: Current  
481        status and future trends. *J. Chromatogr. B* **2008**, *871*, 143 – 163, Hyphenated Tech-  
482        niques for Global Metabolite Profiling.
- 483 (12) Samanipour, S.; Reid, M. J.; Bæk, K.; Thomas, K. V. Combining a Deconvolution and  
484        a Universal Library Search Algorithm for the Nontarget Analysis of Data-Independent  
485        Acquisition Mode Liquid Chromatography-High-Resolution Mass Spectrometry Re-  
486        sults. *Environ. Sci. Technol.* **2018**, *52*, 4694–4701.
- 487 (13) Samanipour, S.; Kaserzon, S.; Vijayasathy, S.; Jiang, H.; Choi, P.; Reid, M. J.;  
488        Mueller, J. F.; Thomas, K. V. Machine learning combined with non-targeted LC-HRMS  
489        analysis for a risk warning system of chemical hazards in drinking water: A proof of  
490        concept. *Talanta* **2019**, *195*, 426 – 432.
- 491 (14) Brack, W.; Hollender, J.; de Alda, M. L.; Müller, C.; Schulze, T.; Schymanski, E.;  
492        Slobodnik, J.; Krauss, M. High-resolution mass spectrometry to complement monitor-  
493        ing and track emerging chemicals and pollution trends in European water resources.  
494        *Environ. Sci. Eur.* **2019**, *31*, 62.
- 495 (15) Minkus, S.; Bieber, S.; Letzel, T. Spotlight on mass spectrometric non-target screening  
496        analysis: Advanced data processing methods recently communicated for extracting,  
497        prioritizing and quantifying features. *Analytical Science Advances* **2022**, *3*, 103–112.

- 498 (16) van Herwerden, D.; O'Brien, J. W.; Choi, P. M.; Thomas, K. V.; Schoenmakers, P. J.;  
499 Samanipour, S. Naive Bayes classification model for isotopologue detection in LC-  
500 HRMS data. *Chemometrics and Intelligent Laboratory Systems* **2022**, *223*, 104515.
- 501 (17) Gertsman, I.; Barshop, B. A. Promises and pitfalls of untargeted metabolomics. *Journal*  
502 *of Inherited Metabolic Disease* **2018**, *41*, 355–366.
- 503 (18) Watson, D. G. A ROUGH GUIDE TO METABOLITE IDENTIFICATION USING  
504 HIGH RESOLUTION LIQUID CHROMATOGRAPHY MASS SPECTROMETRY  
505 IN METABOLOMIC PROFILING IN METAZOANS. *Computational and Structural*  
506 *Biotechnology Journal* **2013**, *4*, e201301005.
- 507 (19) Fedorova, E. S.; Matyushin, D. D.; Plyushchenko, I. V.; Stavrianidi, A. N.;  
508 Buryak, A. K. Deep learning for retention time prediction in reversed-phase liquid  
509 chromatography. *Journal of Chromatography A* **2022**, *1664*.
- 510 (20) Rigano, F.; Arigò, A.; Oteri, M.; Tella, R. L.; Dugo, P.; Mondello, L. The retention  
511 index approach in liquid chromatography: An historical review and recent advances.  
512 *Journal of Chromatography A* **2021**, *1640*.
- 513 (21) Smith, R. M. *Chapter 3 Retention index scales used in high-performance liquid chro-*  
514 *matography*; 1995; pp 93–144.
- 515 (22) Aalizadeh, R.; Thomaidis, N. S.; Bletsou, A. A.; Gago-Ferrero, P. Quantitative Struc-  
516 ture–Retention Relationship Models To Support Nontarget High-Resolution Mass Spec-  
517 trometric Screening of Emerging Contaminants in Environmental Samples. *Journal of*  
518 *Chemical Information and Modeling* **2016**, *56*, 1384–1398.
- 519 (23) Lamparczyk, H.; Radecki, A. The role of electric interactions in the retention index con-  
520 cept; Implications in quantitative structure-retention studies. *Chromatographia* **1984**,  
521 *18*, 615–618.



- 522 (24) Farkas, O.; Zenkevich, I. G.; Stout, F.; Kalivas, J. H.; Héberger, K. Prediction of reten-  
523 tion indices for identification of fatty acid methyl esters. *Journal of Chromatography A*  
524 **2008**, *1198-1199*, 188–195.
- 525 (25) Samanipour, S.; O'Brien, J. W.; Reid, M. J.; Thomas, K. V.; Praetorius, A. From  
526 Molecular Descriptors to Intrinsic Fish Toxicity of Chemicals: An Alternative Approach  
527 to Chemical Prioritization. *Environmental Science and Technology* **2022**,
- 528 (26) Boelrijk, J.; van Herwerden, D.; Ensing, B.; Forré, P.; Samanipour, S. Predicting RP-  
529 LC retention indices of structurally unknown chemicals from mass spectrometry data.  
530 *Journal of Cheminformatics* **2023**, *15*, 1–12.
- 531 (27) Alygizakis, N.; Konstantakos, V.; Bouziotopoulos, G.; Kormentzas, E.; Slobodnik, J.;  
532 Thomaidis, N. S. A Multi-Label Classifier for Predicting the Most Appropriate Instru-  
533 mental Method for the Analysis of Contaminants of Emerging Concern. *Metabolites*  
534 **2022**, *12*.
- 535 (28) Hall, L. M.; Hill, D. W.; Menikarachchi, L. C.; Chen, M. H.; Hall, L. H.; Grant, D. F.  
536 Optimizing artificial neural network models for metabolomics and systems biology: An  
537 example using HPLC retention index data. *Bioanalysis* **2015**, *7*, 939–955.
- 538 (29) Aalizadeh, R. et al. Development and Application of Liquid Chromatographic Retention  
539 Time Indices in HRMS-Based Suspect and Nontarget Screening. *Analytical Chemistry*  
540 **2021**, *93*, 11601–11611.
- 541 (30) Aalizadeh, R.; Nikolopoulou, V.; Thomaidis, N. S. Development of Liquid Chromato-  
542 graphic Retention Index Based on Cocamide Diethanolamine Homologous Series (C(n)-  
543 DEA). *Analytical Chemistry* **2022**, *94*, 15987–15996.
- 544 (31) van Herwerden, D.; Samanipour, S. Dataset for: RPLC chemical space  
545 prediction. 2023; [https://figshare.com/articles/dataset/Dataset\\_for\\_RPLC\\_](https://figshare.com/articles/dataset/Dataset_for_RPLC_chemical_space_prediction/22587262)  
546 [chemical\\_space\\_prediction/22587262](https://figshare.com/articles/dataset/Dataset_for_RPLC_chemical_space_prediction/22587262).

- 547 (32) PubChem substructure fingerprint. [https://ftp.ncbi.nlm.nih.gov/pubchem/](https://ftp.ncbi.nlm.nih.gov/pubchem/specifications/pubchem_fingerprints.pdf)  
548 [specifications/pubchem\\_fingerprints.pdf](https://ftp.ncbi.nlm.nih.gov/pubchem/specifications/pubchem_fingerprints.pdf).
- 549 (33) Yap, C. W. PaDEL-descriptor: An open source software to calculate molecular descrip-  
550 tors and fingerprints. *Journal of Computational Chemistry* **2011**, *32*, 1466–1474.
- 551 (34) Yang, F.; van Herwerden, D.; Preud'homme, H.; Samanipour, S. Collision Cross Section  
552 Prediction with Molecular Fingerprint Using Machine Learning. *Molecules* **2022**, *27*.
- 553 (35) Barron, L.; Loftus, N. Gradient retention time predictopm for 653 pesticides on a  
554 biphenyl column using machine learning. *Chromatography Today* **2019**,
- 555 (36) Klamt, A.; Eckert, F.; Reinisch, J.; Wichmann, K. Prediction of cyclohexane-water dis-  
556 tribution coefficients with COSMO-RS on the SAMPL5 data set. *Journal of Computer-*  
557 *Aided Molecular Design* **2016**, *30*.
- 558 (37) Aalizadeh, R.; von der Ohe, P. C.; Thomaidis, N. S. Prediction of acute toxicity of  
559 emerging contaminants on the water flea *Daphnia magna* by Ant Colony Optimization-  
560 Support Vector Machine QSTR models. *Environmental science. Processes impacts*  
561 **2017**, *19*, 438–448.

562 **TOC Graphic**

563

