

# 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, 1098 XH, the Netherlands*

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

<sup>¶</sup>*Queensland Alliance for Environmental Health Sciences (QAEHS), The University of Queensland, Brisbane, QLD 4102, Australia*

<sup>§</sup>*UvA Data Science Center, University of Amsterdam, Amsterdam, 1012 WP, the Netherlands*

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

## Abstract

1  
2 The chemical space is comprised of a vast number of possible structures, of which  
3 an unknown portion comprises the human and environmental exposome. Such sam-  
4 ples are frequently analyzed using non-targeted analysis via liquid chromatography  
5 (LC) coupled to high-resolution mass spectrometry often employing a reversed phase  
6 (RP) column. However, prior to analysis, the contents of these samples are unknown  
7 and could be comprised of thousands of known and unknown chemical constituents.  
8 Moreover, it is unknown which part of the chemical space is sufficiently retained and

9 eluted using RPLC. Therefore, we present a generic framework that uses a data driven  
10 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>

36

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, meaning that these  
97 compounds can generally be difficult to analyze and require specific methods. All chemi-  
98 cals in between the ‘maybe’ regions are classified as ‘inside’. For the RPLC classification  
99 model, a dataset with chemicals that were ‘inside’, ‘maybe’ inside, and ‘outside’ of the RPLC  
100 subspace was constructed (Figure 1B). Finally, the application of the RPLC classification  
101 model was showcased by applying it on the NORMAN SusDat database, which is a collec-  
102 tion of expert curated environmentally relevant chemicals that have been actively used for  
103 screening of complex samples. All training and test datasets for constructing the models and  
104 the NORMAN SusDat database with the calculated fingerprints can be found on Figshare.<sup>31</sup>

105

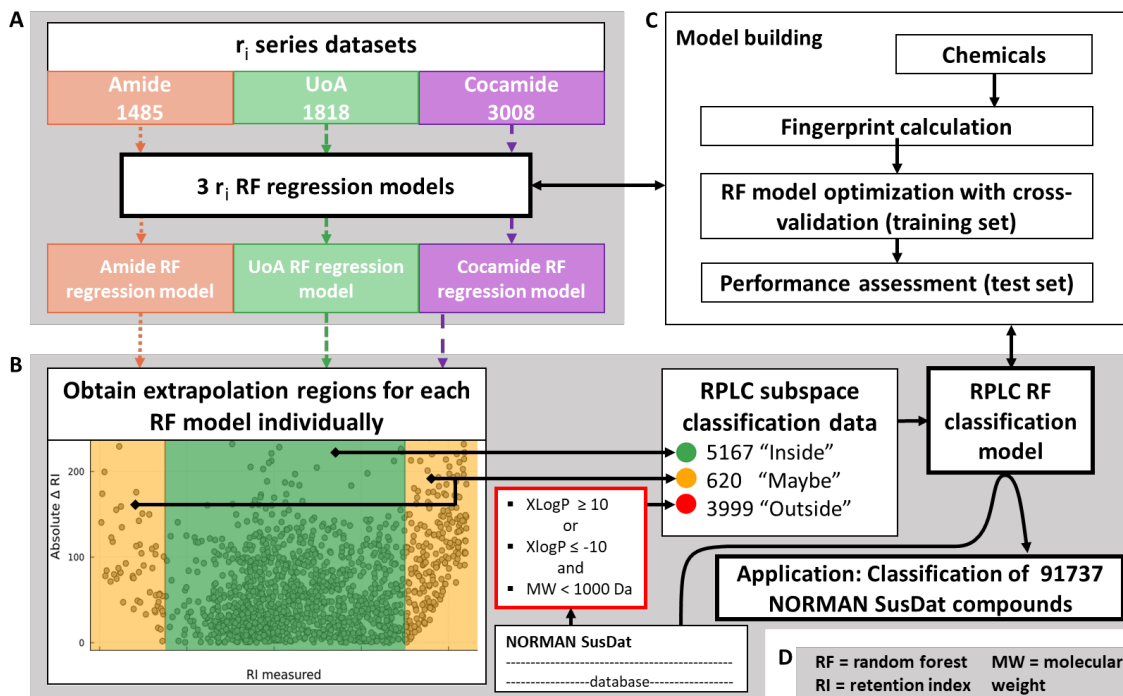


Figure 1: Workflow for construction of the RPLC classification model, comprising of the construction of three  $r_i$  RF regression models (A, Section ‘Retention Index Random Forest Regression Models’) and the construction the RPLC dataset for the RPLC RF classification model, which was applied to NORMAN SusDat (B, Section ‘RPLC Random Forest Classifier’ and ‘RPLC Space Prediction for NORMAN SusDat’). C shows the model setup (Section ‘Fingerprint Calculations’ and ‘Retention Index Random Forest Regression Models’) and D contains an overview of the abbreviations.

## 106 Fingerprint Calculations

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

116 used (i.e., PubChem fingerprint 115 - 262). These fingerprints were converted and reduced  
117 to a total of 10 additional variables, which were the number of rings with a size of 3, 4,  
118 5, 6, 7, 8, 9, 10, the number of aromatic rings, and the number of hetero-aromatic rings.  
119 Since the PubChem fingerprints are binary, there were multiple columns describing the same  
120 information but only differing in the number of a ring of a certain size. For example, for a  
121 ring size of 3, there were 2 fingerprints, namely PubChem fingerprint 115 and 122, which  
122 were described as more than 1 ring with a size of 3 or more than 2 rings with a size of 3,  
123 respectively. In case a molecule contained 2 rings with a size of 3, the PubChem fingerprints  
124 115 would be 0 and 122 would be 1, which was converted to a single variable for our model  
125 containing the number of rings with a size of 3, meaning that this variable would be equal to  
126 2 for this example case. An overview of which PubChem fingerprints were used for each of  
127 the 10 reduced PubChem variables can be found in table S2. Finally, it should be noted that  
128 the use of canonical SMILES for these type of fingerprints would yield no different result  
129 compared to stereoisomeric SMILES, as atom distances and number of rings will remain  
130 consistent.

131

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

133 To show that fingerprints can be used to describe retention behavior in RPLC and for set-  
134 ting up the dataset for the RPLC classification model, random forest (RF) regression models  
135 were built using three different retention index series (Figure 1A). The three series used for  
136 this, were the amide<sup>28</sup>, University of Athens (UoA)<sup>29</sup>, and cocamide series.<sup>30</sup> For each of  
137 the series, the measured  $r_i$  were obtained from their respective articles, yielding 1485, 1818,  
138 and 3008 unique chemicals with measured  $r_i$  values for the amide, UoA, and cocamide series,  
139 respectively. For all chemicals, the 2DAPC and PubChem fingerprints were calculated ac-  
140 cording to Section ‘Fingerprint Calculations’. For each  $r_i$  series, data was split into a training  
141 and test set, at random, with a ratio of 0.85:0.15, ensuring similar coverage of the  $r_i$  range

142 in both sets. The test set was only used for testing and thus never used for training. For  
143 optimization of the RF regression models, the training set was used with a 0.8:0.2 split for  
144 training and cross-validation, respectively. This ratio of split has been shown to be effective  
145 in such data sets.<sup>25,26,34,35</sup> The RF regression models used a third of the features (i.e., 264)  
146 for training each tree. The parameters that were optimized were the minimum number of  
147 samples per leaf and the number of trees. The minimum number of samples per leaf tested  
148 were 4, 6, 8, 10, 15, and 20. The tested number of trees were 50, 100, 150, 200, 250, 300,  
149 350, 400, 500, 600, 700, 800, 900, and 1000. In addition, the random state for splitting the  
150 cross-validation set and selection of the features in the RF models for each tree was also  
151 varied with values of 1, 2, and 3. The accuracy of the cross-validation set for each possible  
152 combination of the minimum number of samples per leaf, number of trees, and random state  
153 was used for the optimization of the RF models. After obtaining the optimized models for  
154 the amide, UoA, and cocamide series, the applicability domains were assessed according to  
155 Section ‘Applicability Domain Calculations’. Finally, for each  $r_i$  series, the optimized model  
156 and applicability domain assessment were applied on the test set to evaluate the performance  
157 of the model on unseen data.

158

## 159 **RPLC Random Forest Classifier**

160 The dataset for building the RPLC classifier model was comprised of three classes: ‘inside’,  
161 ‘maybe’, and ‘outside’ the RPLC subspace (Figure 1B). The ‘outside’ chemicals were ob-  
162 tained from the NORMAN SusDat database based on their extreme XLogP values, assuming  
163 that these cannot be analysed using RPLC regardless of the method used. Here, the XLogP  
164 was chosen rather than the logD due to the fact that it is easier to predict, more stable,  
165 and more accurate.<sup>36</sup> For the ‘outside’ case, a total of 3999 compounds with a XLogP value  
166 above 10 or below -10 and with a molecular weight below 1000 Da were obtained. As for  
167 the ‘inside’ and ‘maybe’ chemicals, these were obtained from the experimentally defined  $r_i$

168 values by the three  $r_i$  series. For each of the series, the absolute difference between the  
169 predicted and measured  $r_i$  (i.e., the residuals) versus the measured  $r_i$  values were plotted  
170 and the regions of extrapolation were identified. These regions were obtained based on the  
171 increasing residuals that were caused by the inherent over estimation and under estimation  
172 of a RF regression model, which are associated with either extremely low or extremely high  
173  $r_i$  values, respectively. These regions correspond to chemicals that elute close to  $t_0$  or are  
174 very difficult to elute from the column (i.e., require a relatively high percentage of organic  
175 modifier). The chemicals with a measured  $r_i$  in these extrapolation regions were labeled  
176 as ‘maybe’ and the remaining chemicals were labeled as ‘inside’ the RPLC subspace. This  
177 yielded a total of 620 ‘maybe’ and 5167 ‘inside’ compounds. Whenever a chemical SMILES  
178 was found in multiple classes (i.e., it was present in multiple datasets of the  $r_i$  models), it  
179 was removed from the lower ranking RPLC classes and kept in the highest ranking RPLC  
180 class (i.e., ‘inside’ > ‘maybe’ > ‘outside’ RPLC class rank). For example, if a chemical was  
181 found in the ‘maybe’ region for UoA and in the ‘inside’ for Cocamide, it would be classified  
182 as ‘inside’. More details on the division between the ‘inside’ and ‘maybe’ classification can be  
183 found in Section ‘RPLC Classification Model’ as these are based on the results of the three  
184 RF regression models. It should be noted, that even though output information from the  $r_i$   
185 models has been used to set up the classification dataset, the regression and classification  
186 models are independent, meaning that there is no data leakage taking place.

187

188 The calculated fingerprints (Section ‘Fingerprint Calculations’) for the dataset described  
189 above were used for building the RPLC classifier model with a training set/test set split of  
190 0.85:0.15, ensuring equal distribution of each class in both sets. The optimized RF classifier  
191 model was obtained using the same approach as for the RF regression models (see Section  
192 ‘Retention Index Random Forest Regression Models’). For this model, the applicability  
193 domain was also obtained as described below. Finally, the optimized RPLC classification  
194 model and applicability domain assessment was applied to the test set and the performance

195 was evaluated.

196

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

198 To showcase the model's potential, it was applied to the NORMAN SusDat database.<sup>5</sup> For  
199 this, the 2DAPC and reduced PubChem fingerprints for a total of 91737 chemicals with a  
200 molecular weight below 1000 Da from SusDat were calculated. These fingerprints were then  
201 used to calculate the leverage of each chemical with the RPLC classifier training set, as  
202 explained in the next section 'Applicability Domain Calculations', and to apply the RPLC  
203 classifier model to each of the SusDat chemicals. To visualize the coverage of each class  
204 (i.e., 'inside', 'maybe', and 'outside' the RPLC subspace), the molecular weight was plotted  
205 against the XLogP, which were obtained from the descriptor calculations of PaDEL.

206

## 207 **Applicability Domain Calculations**

208 Applicability domain calculations were used to assess whether the training data, used in the  
209 random forest models, sufficiently covered the variable space for new chemicals on which the  
210 models need to be applied.<sup>25,37</sup> This was done through leverage calculations of a chemical  
211 with the entire training set, yielding a distance of that chemical to the training set. Finger-  
212 prints are used to calculate this distance, meaning that lower distance values are obtained for  
213 compound that are structurally more similar to the training set than compounds with high  
214 leverage values. Equation 1 shows how the leverage is calculated, where  $X$  is the training  
215 data matrix and  $x_i$  is the sample vector, both containing the 2DAPC and reduced PubChem  
216 fingerprints for our models. To set a threshold for this, the leverage was calculated for all  
217 training samples with the entire training set of a model, yielding values between 0 and 1.  
218 Then, a leverage threshold was obtained that covered 95% of the training data. If a chemical,  
219 compared to the training set of the model in question, had a value lower than the leverage

220 threshold, the compound was within the applicability domain, and, if the value was above  
221 the leverage threshold, the results should be taken with care as the training data might not  
222 be sufficiently describing the variable space for the new compound.

223

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

## 224 **Calculations and Code Availability**

225 The calculations and development of the models were executed on a personal computer with  
226 12 CPUs and 32 GB of RAM, using Windows 10. The  $r_i$  regression and RPLC classifi-  
227 cation models were developed and evaluated with the Julia programming language (v1.6).  
228 The code for using the  $r_i$  regression models and RPLC space prediction model is available  
229 at: [https://bitbucket.org/Denice\\_van\\_Herwerden/riprediction/src/main/](https://bitbucket.org/Denice_van_Herwerden/riprediction/src/main/). This Ju-  
230 lia package contains functions for obtaining the required 2DAPC and reduced PubChem  
231 fingerprints and for using the  $r_i$  regression models and RPLC sub space classification model.

232

## 233 **Results and discussion**

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

235 All three  $r_i$  regression models obtained an accuracy of 81% for the training set and, for the  
236 test set, the amide, UoA, and cocamide models had an accuracy of 68%, 70%, and 67%,  
237 respectively. The  $r_i$  regression models were built and optimized for the amide, UoA and  
238 cocamide series. Grid optimization of each of these models showed that the number of trees  
239 did not influence the performance of the model (Figures S1, S2, and S3). Therefore, to  
240 keep the model light, 200 trees were selected. As for the minimum number of samples per  
241 leaf, 8 was found to be the optimum, based on the training and cross-validation accuracy.

242 When evaluating the predicted versus the measured  $r_i$  values for these models a trend of over  
243 prediction for lower  $r_i$  values and under prediction of higher  $r_i$  values was found (Figures S4,  
244 S6, and S8), corresponding to the regions where the RF regression models were extrapolat-  
245 ing. These regions were used for establishing the ‘maybe’ areas for the RPLC classification  
246 dataset.

247

248 Most compounds (i.e., 88.5%) in our test set appeared to be within the applicability  
249 domain of each model. To obtain the applicability domains of these models, a 95% leverage  
250 threshold of 0.189 for amide, 0.652 for UoA, and 0.424 for cocamide was found for the train-  
251 ing sets. For the training set the leverage values range between 0 and 1, meaning that the  
252 lower threshold for the amide model showed how similar most of the amide compounds were  
253 to each other, while for the UoA and cocamide models, the higher thresholds corresponded  
254 with the larger variety of chemical structures found in the dataset. When the leverage cal-  
255 culations were applied on the test sets for these models, a total of 22, 34, and 54 compounds  
256 were found to be outside of the applicability domain for the amide, UoA, and cocamide  $r_i$   
257 models, respectively. This does not necessarily mean that the predicted outcome for these  
258 cases was wrong, as can be seen in figures S4, S6, and S8. Here, most chemicals outside the  
259 applicability domain still follow the trend of the other data points. However, the outcome  
260 should be taken with care as the model might insufficiently cover the chemical space for a  
261 new compound in question, especially for leverage values  $> 1$ . It should be noted that the  
262 largest training set leverage value obtained from our applicability domain calculations was 1.

263

264 The cocamide RF regression model used the most fingerprints for the prediction of the  
265  $r_i$  indices (i.e., 215 fingerprints), while the UoA and amide  $r_i$  models used 165 and 61, re-  
266 spectively. The low number of fingerprints used for amide was not surprising due to the  
267 fact that the compounds in this  $r_i$  series are only comprised of C, H, N, and O. Hence, the  
268 amide  $r_i$  model only used the 2DAPC fingerprint counts with a certain distance between C,



269 N, and O atoms. At first sight, this was also noticeable when comparing the top 20 most  
270 important fingerprints for the three  $r_i$  models (S3). The most contributing fingerprints for  
271 the amide  $r_i$  model were the distances 1 till 7 between two C atoms with importance ranging  
272 between 27% and 4%. As for the UoA  $r_i$  model, C-Cl and C-X distance begin to contribute  
273 more to the model and the most important fingerprint (i.e., distance 7 between C-C) only  
274 contributes 9.6%, having an overall more divided importance between a larger group of con-  
275 tributing features than the amide model. Finally, a similar trend was also observed for the  
276 cocamide model, except that the C-X distances start to play a more important role than the  
277 C-Cl distances, which could be explained by the higher number of halogens present in the  
278 compounds from the cocamide dataset. This variability in important features used in each  
279  $r_i$  regression model shows that different structures may be better captured by one  $r_i$  model  
280 vs another, due to the diversity of training set in terms of chemical structures. This, also,  
281 further indicates the need for a more generic model incorporating the information from all  
282 three  $r_i$  models.

283

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

289

## 290 **RPLC Classification Model**

291 To build the RPLC classification model, it was assumed that the chemicals are in solution  
292 and that the chemicals can be injected into a system. Additionally, the model focuses on  
293 whether an analyte could be analyzed with RPLC regardless of experimental parameters  
294 or sample pretreatment. The dataset for this was comprised of 5167 ‘inside’, 620 ‘maybe’

295 inside, and 3999 ‘outside’ chemicals for the RPLC subspaces. The ‘outside’ cases were ob-  
296 tained from NORMAN SusDat with extreme XLogP values, while the ‘inside’ and ‘maybe’  
297 cases came from the three  $r_i$  regression models. In figures S10, S11, and S12 the extrapola-  
298 tion limits for each of the models are defined. For  $r_i$  range for the ‘inside’ RPLC subspace  
299 for the amide, UoA, and cocamide series were 350-900, 100-900, and 250-1300, respectively.  
300 Each of the  $r_i$  series has their own scale and range of retention index values. Therefore, these  
301 values are not directly comparable between the series. All compounds that had a higher or  
302 lower  $r_i$  value for the corresponding range of the model it was coming from, were classified  
303 as ‘maybe’ inside the RPLC subspace, due to the fact that these chemicals either elute close  
304 to  $t_0$  or require high percentages of organic eluent to be eluted.

305

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

319

320 As for the applicability domain of the RPLC classification model, the 95% leverage  
321 threshold of the training set was 0.209 (Figure S14). In total, 102 compounds from the test

322 set (i.e., 6.9%) had a leverage with the training set that was higher than 0.209, of which 31  
323 cases had leverage values above 1. Out of these 102 cases only 10 were wrongly classified  
324 and had leverage values ranging between 0.209 to the most extreme (i.e., 809.255), showing  
325 that in this case higher leverage values did not necessarily mean that the model would have  
326 a higher error. However, it should be noted that cases with a very large leverage should be  
327 considered with extra care, as they may have a higher level of uncertainty.

328

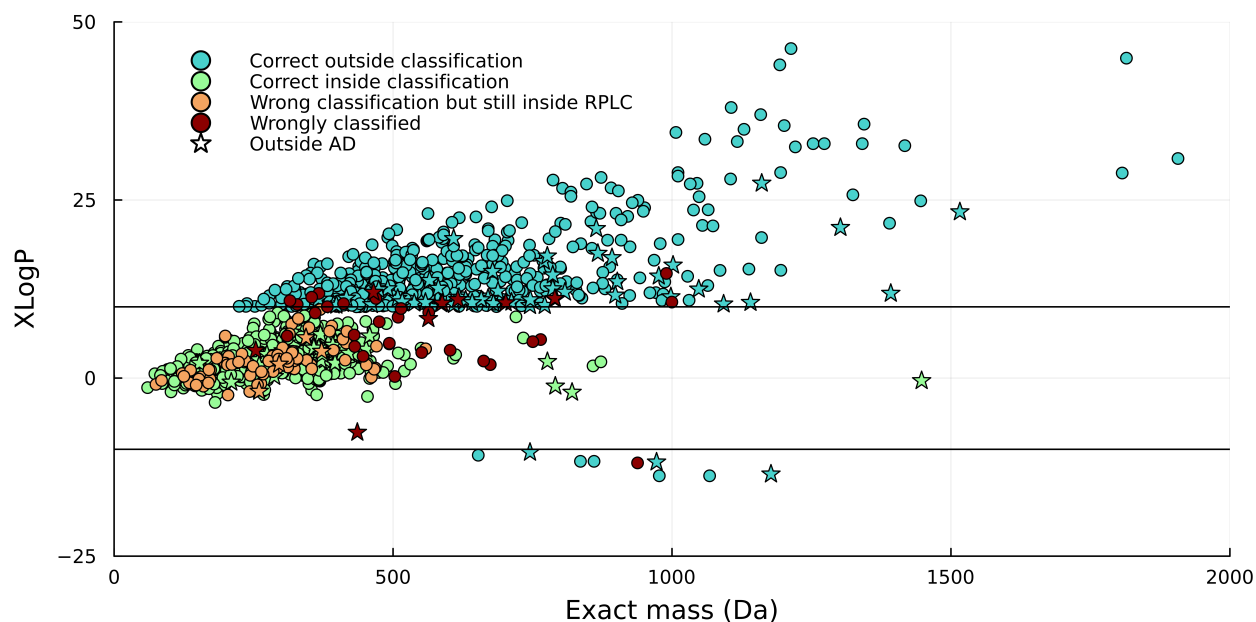


Figure 2: XLogP values versus the molecular weight for the RPLC classification test set. In blue are the correctly classified ‘outside’ cases, in green are the correctly classified ‘inside’ and ‘maybe’ cases, in orange are the wrongly classified ‘inside’ cases as ‘maybe’ and vice versa, in red the wrongly 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

329 A total of 280 features were contributing to the RPLC classification model. This is more  
330 than for each of the three  $r_i$  regression models, which was expected due to the higher variety  
331 in chemical structures used in the RPLC classification model. The 20 most contributing  
332 features are mainly described by ring related features and distances between combinations  
333 C, N, and O atoms. A previous version of the model that was tested, using only the 2DAPC

334 fingerprints, frequently wrongly classified ‘inside’ as ‘outside’ due to the high degree of cyclic-  
335 ity in the chemical structures (e.g., peimine). Hence, the addition of the reduced PubChem  
336 fingerprints better captures these chemical properties. As a result, the number of rings with  
337 a size of 6, the minimum number of aromatic rings, and the number of rings with a size of  
338 5 were also part of the top 20 most contributing features.

339

340 In total, considering the extreme misclassifications, 9 out of 599 ‘outside’ chemicals were  
341 wrongly classified as ‘inside’ or ‘maybe’ inside the RPLC subspace and 14 out of the 767  
342 ‘inside’ and 12 out of the 102 ‘maybe’ cases were classified as ‘outside’ the RPLC subspace.  
343 Two of the nine wrongly classified ‘outside’ cases were organic complexes that, in the mobile  
344 phase, would be analyzed as multiple smaller molecules (e.g., Gadopentetic acid dimeglu-  
345 mine salt). Also, another case was a surfactant containing a positive and negative charge  
346 (i.e., 4-Dodecyl-2-[(2-nitrophenyl)azo]phenol). This case was a chemical that falls ‘outside’  
347 of the RPLC space due to its predicted XLogP value of 10.452. However, the charges on  
348 this molecule would make it difficult to calculate this value accurately. Lexidronam was one  
349 of the ‘maybe’ cases that was classified as ‘outside’, due to a large leverage value of 26.0  
350 and the fact that it elutes at  $t_0$  (i.e., amide scale  $r_i$  of 206 versus urea  $r_i = 200$ ), indicating  
351 the need for special gradients to be able to retain such a chemical. As for the ‘inside’ cases  
352 that were wrongly classified as ‘outside’, generally larger, branched (e.g., SCHEMBL312614),  
353 or hydrolyzing (e.g., Bis[2-(perfluorohexyl)ethyl] Phosphate, respectively) chemicals showed  
354 higher likelihood of such misclassifications. Again these are structures that may require very  
355 specific adjustment of experimental condition (e.g., pH of mobile phase) to fit them within  
356 the RPLC analyzable chemical subspace.

357

358 Overall, our RPLC classification model was highly successful in identifying the chemical  
359 structures that are easily analyzable via RPLC (i.e., ‘inside’ cases) as well as the ‘maybe’ and  
360 ‘outside’ cases. The classification model used a combination of similar molecular fingerprints

361 as those used by the three  $r_i$  models, taking advantage of all the structural information.

## 362 **NORMAN SusDat Chemical Space Prediction**

363 Finally, the RPLC classification model was applied to a set of small molecules (i.e., molecular  
364 weight < 1000) from the NORMAN SusDat database. In total, 80503 chemicals were within  
365 the applicability domain with leverage values  $\leq 0.209$ , 6570 compounds had leverage values  
366 between 0.209 and 1, and 4664 compounds had even larger leverages. This showed that the  
367 RPLC classification model was suitable for a large variety, 87.8%, of compounds present in  
368 SusDat. The model predicted that 79.0% of the compounds would fit ‘inside’ the RPLC  
369 subspace, 2.0% was ‘maybe’ in this space, and 19.1% was ‘outside’ of the RPLC subspace.  
370 Examples of molecules classified as ‘inside’, ‘maybe’, and ‘outside’ were carbamazepine, su-  
371 dan I, and coronene, respectively. When comparing the relationship between XlogP and  
372  $r_i$ , it is clearly observable that these parameters, even though relatively linear, are insuffi-  
373 cient to determine if a chemical fits the RPLC subspace, figure 3. In figures S16,S17, and  
374 S18, the XlogP values of the chemicals with the same  $r_i$  range vary between -10 to +10 units.

375

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

382

383 The amide  $r_i$  model is the least suited scale based on its applicability domain coverage  
384 since only 44500 (i.e., 48.5%) chemicals fell within the applicability domain. For the chem-  
385 icals that were outside the applicability domain, 18988 had a leverage value between 0.189  
386 and 1 (i.e., similar to the full training set) and 28249 had an even higher leverage value. As

387 for the UoA and cocamide  $r_i$  models, 71022 (i.e., 77.4%) and 74252 (i.e., 80.9%) compounds  
388 were within the applicability domain. For the UoA model, 3421 and 17294 chemicals had a  
389 leverage value below and above 1, respectively, and the cocamide model had 5947 chemicals  
390 with a leverage value below 1 and 11538 chemicals with higher leverage values. Figures S16,  
391 S17, and S18 show the coverage of the ‘inside’, ‘maybe’, and ‘outside’ RPLC classes in terms  
392 of the XLogP values versus the predicted  $r_i$  values for the amide, UoA, and cocamide series.  
393 As expected the chemicals classified as ‘maybe’ inside RPLC are mainly clustering around  
394 the lower and higher  $r_i$  values. While the chemicals classified as ‘outside’ the RPLC space  
395 span the entire  $r_i$  range for each of the three  $r_i$  series, suggesting that  $r_i$  prediction would  
396 also be insufficient to define the boundaries of the RPLC subspace.

397

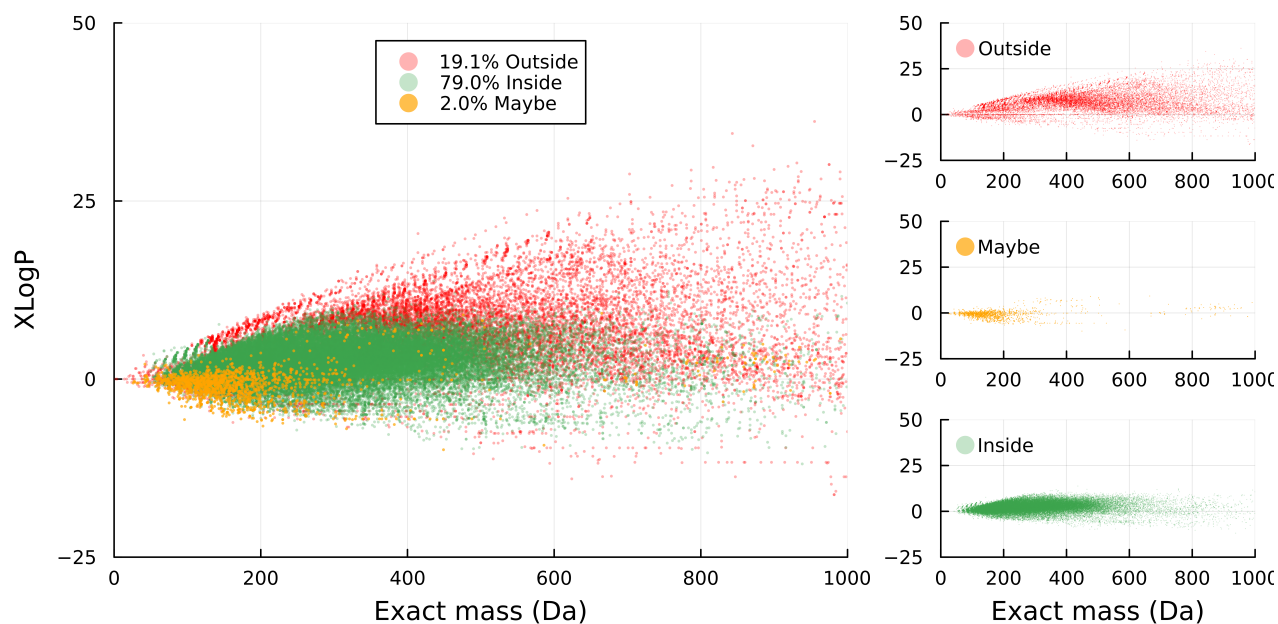


Figure 3: XLogP values versus the molecular weight for the NORMAN SusDat database compounds with a molecular weight below 1000 Da. In red, orange, and green are the compounds that were classified as ‘outside’, ‘maybe’, and ‘inside’ the RPLC chemical space, respectively. The subplots on the left show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

## Potentials and Limitations

Overall, we developed four models for exploration of the RPLC subspace. The  $r_i$  regression models showed that fingerprints can be used for describing RPLC retention indices. Consequently, these fingerprints were used for RPLC classification model building. This model was able to predict whether chemicals were ‘inside’, ‘maybe’ inside, or ‘outside’ of RPLC chemical subspace with an accuracy of 92% on the test set. Applying the RPLC classification model on NORMAN SusDat showed that 19.1% of the compounds were classified as ‘outside’ the RPLC subspace. This means that, when performing identification on NTA RPLC samples, candidates classified as ‘outside’ compounds are unlikely to be the true structure of the chemical and can be removed to reduce the number of false positive identifications. In terms of suspect screening, it can save computational time since the ‘outside’ chemicals fall ‘outside’ of the RPLC subspace and thus should not be screened for. Additionally, 87.8% of NORMAN SusDat was within the applicability domain of the RPLC classifier, showing good coverage of a variety of compounds. The RPLC classification model also showed that the XLogP or  $r_i$  values alone are not sufficient to define the RPLC subspace.

The RPLC classification model was built with a focus on small organic molecules (i.e.,  $\leq 1000$  Da). The model did overall have more difficulties with regard to more bulky and branched or surfactant-like chemicals as well as metal-organic compounds. Additionally, the model was not able to properly predict the RPLC subspace class of chemicals that are organic complexes, due to the fact that in solution those are dissociated into multiple individual structures. The latter is not a major limitation for the model itself, since, using expert knowledge, they can be easily identified. Generally, as knowledge on analyzable chemicals with RPLC grows, the model could easily be rebuilt and expanded for the range of analytes. In the near future, we are planning to expand our model to other selectivities, such as HILIC, taking advantage of public retention repositories, such as RepoRT.<sup>38</sup> This allows for further understanding of what part of the chemical space is actually covered by the selectivities used

425 in NTA and what we are missing.

426

427 Moreover, the RPLC classification model uses a data driven approach and is intended  
428 for quick screening of the RPLC chemical space. The model assumes that compounds are  
429 analyzable with RPLC regardless of the chemicals solubility, experimental parameters, or  
430 pretreatment steps taken. This means that it cannot be assumed that chemicals ‘inside‘ the  
431 RPLC space will be analyzable with every RPLC method. Here, the method subspace plays  
432 a major role when looking at what individual NTA methods can cover, becoming an even  
433 more complex issue due to the fact that sample pretreatment, gradient program’s, and RP  
434 column selectivities play a large influence on this. Defining the method chemical space would  
435 be the next step in understanding what part of the vast chemical space we are covering and,  
436 more importantly, excluding with our current NTA methods.

437

## 438 **Acknowledgement**

439 The authors thank the Environmental Monitoring and Computational Mass Spectrometry  
440 ([www.emcms.info](http://www.emcms.info)) group for their insights and feedback. The Queensland Alliance for En-  
441 vironmental Health Sciences. Finally, the University of Queensland gratefully acknowledges  
442 the financial support from the Queensland Department of Health. J.W.O is the recipient of  
443 an NHMRC Emerging Leadership Fellowship (EL1 2009209).

## 444 **Supporting Information Available**

445 Overview of performance for using different types of molecular fingerprints, composition of  
446 reduced PubChem fingerprints, optimization, prediction, leverage, and feature importance  
447 results for the 3 RF regression models and the RPLC classification model, and the RPLC  
448 classification of NORMAN SusDat visualized by plotting the XLogP values versus the pre-



449 dicted  $r_i$  values for the three  $r_i$  regression models.

## 450 **Author Information**

451 Corresponding Author:

452 Saer Samanipour

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

454 University of Amsterdam,

455 the Netherlands

456 Email: s.samanipour@uva.nl

457

458 Denice van Herwerden

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

460 University of Amsterdam,

461 the Netherlands

462 Email: d.vanherwerden@uva.nl

463

## 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. *Journal of Cheminformatics* **2017**, *9*, Article number: 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. *Advanced Drug Delivery Reviews* **1997**, *23*, 3–25.
- (7) Bohacek, R. S.; McMartin, C.; Guida, W. C. The art and practice of structure-based drug design: A molecular modeling perspective. *Medicinal Research Reviews* **1996**, *16*, 3–50.
- (8) Hulleman, T.; Turkina, V.; O'Brien, J. W.; Chojnacka, A.; Thomas, K. V.; Samanipour, S. Critical Assessment of the Chemical Space Covered by LC-HRMS Non-Targeted Analysis. *Environmental Science and Technology* **2023**, *57*, 14101–14112.
- (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

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

- 512 (16) van Herwerden, D.; O'Brien, J. W.; Choi, P. M.; Thomas, K. V.; Schoenmakers, P. J.;  
513 Samanipour, S. Naive Bayes classification model for isotopologue detection in LC-  
514 HRMS data. *Chemometrics and Intelligent Laboratory Systems* **2022**, *223*, 104515.
- 515 (17) Gertsman, I.; Barshop, B. A. Promises and pitfalls of untargeted metabolomics. *Journal*  
516 *of Inherited Metabolic Disease* **2018**, *41*, 355–366.
- 517 (18) Watson, D. G. A Rough Guide to Metabolite Identification using High Resolution  
518 Liquid Chromatography Mass Spectrometry in Metabolomic Profiling in Metazoans.  
519 *Computational and Structural Biotechnology Journal* **2013**, *4*, e201301005.
- 520 (19) Fedorova, E. S.; Matyushin, D. D.; Plyushchenko, I. V.; Stavrianidi, A. N.;  
521 Buryak, A. K. Deep learning for retention time prediction in reversed-phase liquid  
522 chromatography. *Journal of Chromatography A* **2022**, *1664*.
- 523 (20) Rigano, F.; Arigò, A.; Oteri, M.; Tella, R. L.; Dugo, P.; Mondello, L. The retention  
524 index approach in liquid chromatography: An historical review and recent advances.  
525 *Journal of Chromatography A* **2021**, *1640*.
- 526 (21) Smith, R. M. *Chapter 3 Retention index scales used in high-performance liquid chro-*  
527 *matography*; 1995; pp 93–144.
- 528 (22) Aalizadeh, R.; Thomaidis, N. S.; Bletsou, A. A.; Gago-Ferrero, P. Quantitative Struc-  
529 ture–Retention Relationship Models To Support Nontarget High-Resolution Mass Spec-  
530 trometric Screening of Emerging Contaminants in Environmental Samples. *Journal of*  
531 *Chemical Information and Modeling* **2016**, *56*, 1384–1398.
- 532 (23) Lamparczyk, H.; Radecki, A. The role of electric interactions in the retention index con-  
533 cept; Implications in quantitative structure-retention studies. *Chromatographia* **1984**,  
534 *18*, 615–618.

- 535 (24) Farkas, O.; Zenkevich, I. G.; Stout, F.; Kalivas, J. H.; Héberger, K. Prediction of reten-  
536 tion indices for identification of fatty acid methyl esters. *Journal of Chromatography A*  
537 **2008**, *1198-1199*, 188–195.
- 538 (25) Samanipour, S.; O'Brien, J. W.; Reid, M. J.; Thomas, K. V.; Praetorius, A. From  
539 Molecular Descriptors to Intrinsic Fish Toxicity of Chemicals: An Alternative Approach  
540 to Chemical Prioritization. *Environmental Science and Technology* **2023**, *57*, 17950–  
541 17958.
- 542 (26) Boelrijk, J.; van Herwerden, D.; Ensing, B.; Forré, P.; Samanipour, S. Predicting RP-  
543 LC retention indices of structurally unknown chemicals from mass spectrometry data.  
544 *Journal of Cheminformatics* **2023**, *15*, 1–12.
- 545 (27) Alygizakis, N.; Konstantakos, V.; Bouziotopoulos, G.; Kormentzas, E.; Slobodnik, J.;  
546 Thomaidis, N. S. A Multi-Label Classifier for Predicting the Most Appropriate Instru-  
547 mental Method for the Analysis of Contaminants of Emerging Concern. *Metabolites*  
548 **2022**, *12*.
- 549 (28) Hall, L. M.; Hill, D. W.; Menikarachchi, L. C.; Chen, M. H.; Hall, L. H.; Grant, D. F.  
550 Optimizing artificial neural network models for metabolomics and systems biology: An  
551 example using HPLC retention index data. *Bioanalysis* **2015**, *7*, 939–955.
- 552 (29) Aalizadeh, R. et al. Development and Application of Liquid Chromatographic Retention  
553 Time Indices in HRMS-Based Suspect and Nontarget Screening. *Analytical Chemistry*  
554 **2021**, *93*, 11601–11611.
- 555 (30) Aalizadeh, R.; Nikolopoulou, V.; Thomaidis, N. S. Development of Liquid Chromato-  
556 graphic Retention Index Based on Cocamide Diethanolamine Homologous Series (C(n)-  
557 DEA). *Analytical Chemistry* **2022**, *94*, 15987–15996.
- 558 (31) van Herwerden, D.; Samanipour, S. Dataset for: RPLC chemical space

- 559 prediction. 2023; [https://figshare.com/articles/dataset/Dataset\\_for\\_RPLC\\_](https://figshare.com/articles/dataset/Dataset_for_RPLC_chemical_space_prediction/22587262)  
560 [chemical\\_space\\_prediction/22587262](https://figshare.com/articles/dataset/Dataset_for_RPLC_chemical_space_prediction/22587262).
- 561 (32) PubChem substructure fingerprint. [https://ftp.ncbi.nlm.nih.gov/pubchem/](https://ftp.ncbi.nlm.nih.gov/pubchem/specifications/pubchem_fingerprints.pdf)  
562 [specifications/pubchem\\_fingerprints.pdf](https://ftp.ncbi.nlm.nih.gov/pubchem/specifications/pubchem_fingerprints.pdf).
- 563 (33) Yap, C. W. PaDEL-descriptor: An open source software to calculate molecular descrip-  
564 tors and fingerprints. *Journal of Computational Chemistry* **2011**, *32*, 1466–1474.
- 565 (34) Yang, F.; van Herwerden, D.; Preud'homme, H.; Samanipour, S. Collision Cross Section  
566 Prediction with Molecular Fingerprint Using Machine Learning. *Molecules* **2022**, *27*.
- 567 (35) Barron, L. P.; Loftus, N. Gradient retention time predictopm for 653 pesticides on a  
568 biphenyl column using machine learning. *Chromatography Today* **2019**,
- 569 (36) Klamt, A.; Eckert, F.; Reinisch, J.; Wichmann, K. Prediction of cyclohexane-water dis-  
570 tribution coefficients with COSMO-RS on the SAMPL5 data set. *Journal of Computer-*  
571 *Aided Molecular Design* **2016**, *30*, 959—967.
- 572 (37) Aalizadeh, R.; von der Ohe, P. C.; Thomaidis, N. S. Prediction of acute toxicity of  
573 emerging contaminants on the water flea *Daphnia magna* by Ant Colony Optimization-  
574 Support Vector Machine QSTR models. *Environmental science. Processes impacts*  
575 **2017**, *19*, 438–448.
- 576 (38) Kretschmer, F.; Harrieder, E.-M.; Hoffmann, M. A.; Böcker, S.; Witting, M. RepoRT:  
577 a comprehensive repository for small molecule retention times. *Nature Methods* **2024**,

578 **TOC Graphic**

579

