

Critical assessment of the chemical space covered by LC-HRMS non-targeted analysis

Tobias Hulleman,^{†,||} Viktoriia Turkina,^{*,†,||} Jake W. O'Brien,^{‡,†} Aleksandra Chojnacka,[†] Kevin V. Thomas,[‡] and Saer Samanipour^{*,†,¶,§}

[†]*Van 't Hoff Institute for Molecular Sciences (HIMS), University of Amsterdam, 1090 GD, Amsterdam, the Netherlands*

[‡]*Queensland Alliance for Environmental Health Sciences (QAEHS), The University of Queensland, 20 Cornwall Street, Woolloongabba, QLD, 4102, Australia*

[¶]*UvA Data Science Center, University of Amsterdam, Amsterdam*

[§]*Queensland Alliance for Environmental Health Sciences (QAEHS), 20 Cornwall Street, Woolloongabba, QLD, 4102, Australia*

|| Contributed equally to this work

E-mail: v.turkina@uva.nl; s.samanipour@uva.nl

Abstract

1
2 Non-targeted analysis (NTA) has emerged as a valuable approach for comprehensive
3 monitoring of chemicals of emerging concern (CECs) in the exposome. The NTA
4 approach, theoretically, is able to identify compounds with diverse physicochemical
5 properties and sources. Non-targeted analysis methods, even though generic and wide
6 scoping, have been shown to have limitations in terms of their coverage of the chemical
7 space, as the number of the identified chemicals in each sample is very low (e.g. $\leq 5\%$).
8 Investigating the chemical space covered by each NTA assay is crucial for understanding
9 the limitations and challenges associated with the workflow from experimental methods

10 to the data acquisition and data processing. In this review, we examined recent NTA
11 studies published between 2017 and 2023 that employed liquid chromatography-high
12 resolution mass spectrometry. The parameters used in each study were documented
13 and reported chemicals at the confidence level 1 and 2 were retrieved. The chosen
14 experimental setups and the quality of reporting were critically evaluated and discussed.
15 The findings revealed that only around 2% of the estimated chemical space was covered
16 by the NTA studies investigated. Little to no trend was found between the experimental
17 setup and the observed coverage, due to the generic and wide scope of NTA studies.
18 The limited coverage of chemical space by the NTA studies highlights the necessity
19 for a more comprehensive approach in experimental and data processing setups to
20 enable the exploration of a broader range of chemical space, with the ultimate goal of
21 protecting human and environmental health. Recommendations to further explore a
22 wider range of the chemical space were given.

23 **Synopsis**

24 The coverage of chemical space via non-target analysis studies and the impact of the exper-
25 imental conditions on that is critically assessed

26 **Introduction**

27 The exposome is the measure of all the exposures, both chemical and non-chemical, of an
28 individual in a lifetime and how those exposures relate to health¹. The chemical space of
29 exposome refers to the chemical space relevant to human and environmental exposure²⁻⁴. On
30 the other hand, the chemical space generally refers to all possible organic structures, that are
31 plausible from organic chemistry point of view^{2,5}. Theoretical estimates of such structures
32 suggests there are around 10^{60} unique structures with molecular weights less than 500 Da^{5,6}.
33 This theoretical chemical space incorporates both known and unknown unknowns^{2,7} and may

34 include structures that can cause adverse effects depending on their exposure levels. In fact,
35 when looking at the known unknowns (i.e. structures recorded in the chemical databases,
36 but not initially known to be present in the sample), several of them have been shown to
37 have adverse effects on environmental and human health⁸⁻¹¹.

38 Chemical prioritization has been one of the main means for dealing with the diversity
39 of chemical space in the human and environmental exposome^{3,12,13}. This consists of ex-
40 ploration of the literature for measured chemicals and their properties/toxicities as well as
41 national/international chemical registries^{14,15}. A combination of predicted properties and
42 toxicity is used to rank chemicals in the databases based on their potential impact on the
43 environment and human health¹⁶. Chemicals with a high potential of such impact are con-
44 sidered as chemicals of emerging concern (CECs)^{17,18}. To facilitate chemical prioritization,
45 several databases consisting of chemical structures, the associated physicochemical proper-
46 ties (both measured and predicted), and their biological activities have been made publicly
47 available (e.g. PubChem,¹⁹ NORMAN Databases,²⁰ and CompTox¹⁴). However, most of
48 these known unknowns remain unmeasured in environmental and biological matrices due
49 to difficulties associated with the inclusion of such a large number of chemicals in routine
50 monitoring programs^{11,13}.

51 Non targeted analysis (NTA) combined with liquid chromatography coupled with high
52 resolution mass spectrometry (LC-HRMS) is considered as one of the most comprehensive
53 methods for the detection and identification of known and unknown unknowns in complex
54 environmental and biological samples^{21,22}. This approach utilizes a generic and wide scope
55 strategy for the sample preparation and analysis to maximize the coverage of the chemical
56 space of the sample^{2,13,21,23-31}. This typically results in very large and complex datasets (e.g.
57 5 GB per sample) that must be pre-processed prior to the identification workflow³¹⁻³³. The
58 NTA data processing workflows include several steps from data conversion to library search
59 and the confidence assessment of the candidate spectra^{2,23,26-29}. Due to the complexity of
60 such datasets and sheer size of the chemical databases, the NTA workflows are not very

61 sensitive and do not result in a high percentage of identified chromatographic features^{34,35}.
62 A more sensitive but less comprehensive data processing alternative is suspect screening
63 where the chemicals of interest are known prior to the data processing workflow. This
64 approach is more sensitive in terms of limits of detection but is unable to detect unknown
65 unknowns^{20,29,36}. These two strategies are commonly employed together for the screening of
66 complex environmental and biological samples²³.

67 The NTA strategy, even though powerful, has not been widely accepted within the reg-
68 ulatory framework due to reproducibility issues^{30,34,37}. Recent studies have indicated that
69 small changes in both experimental (e.g. data dependent vs data independent acquisition)
70 and data processing parameters may result in different outcomes and thus conclusions^{34,35}.
71 Additionally, a recent study has postulated the potential impact of different experimental
72 parameters on the measured chemical space³⁸. In fact, the aforementioned issues with NTA
73 assays have sparked a debate in the scientific community and have given start to a new wave
74 of data processing tools development^{25,39,40}. Additionally, several efforts have been put into
75 better defining the much needed quality control and assurance for such experiments to be
76 successful in detection and identification of the known and unknown unknowns in complex
77 environmental samples, thus better understanding the coverage of the analyzed chemical
78 space^{23,38,41–44}.

79 Several recently published reviews discuss in detail the impact of different steps on the
80 chemical space coverage through different experimental approaches^{2,23,26,27}. They cover both
81 data processing and experimental parameters including study scope, sampling and sample
82 treatment, instrumental conditions, data processing and treatment, and reporting. However,
83 none of these reviews attempted to assess (i.e. quantify) the coverage of the identified
84 chemical space reached by the already conducted NTA environmental studies. Quantification
85 of the coverage of chemical space by an analytical method is not a trivial task. Theoretically,
86 it can be quantified as the number of identified compounds in the given sample divided
87 by the number of all compounds present in the chemical subspace of the sample. But

88 practically, this calculation is impossible, due to the complex chemical nature of samples
89 and the number of unknown constituents. Nevertheless, the investigation of experimentally
90 explored chemical space is highly relevant for the researchers to be aware of the limited
91 coverage of the associated chemical space.

92 In this review, we aim to quantify the coverage of the identified chemical space by recent
93 environmental studies and investigate the relationship between the selected experimental
94 parameters and the explored chemical space. To quantify the covered chemical space via
95 NTA, we collected all recent studies that perform NTA (not suspect screening) and re-
96 ported levels 1 and 2⁴⁵, in terms of identification, structures. Additionally, we limited the
97 scope of this study to semi-polar and polar chemicals analyzable with liquid chromatogra-
98 phy coupled with high resolution mass spectrometry , resulting in a total of 57 papers.
99 As an approximation of the chemical space the NORMAN SusDat database containing
100 around 60k unique chemicals with available PubChem CIDs (compound ID number) was
101 used (<https://www.NORMAN-network.com/nds/susdat/susdatSearchShow.php>). We
102 collected a list of experimental and instrumental parameters, including sample preparation
103 (i.e. storage and extraction conditions), chromatographic separation (e.g. eluents, gradient
104 type, and injection volume), high resolution mass spectrometry settings (e.g. mass analyzer,
105 data acquisition mode, and polarity), and data processing workflows (e.g. mass and reten-
106 tion time tolerance, retention time domain alignment and databases used for the search).
107 We also noted any unreported parameters to identify the most commonly omitted settings.
108 Furthermore, we extracted information on the scope of the studies and samples analyzed.

109 Finally, we estimated the coverage of chemical space explored by recent NTA studies
110 by comparing the structures identified in these studies with the chemical space represented
111 by the compounds in the NORMAN SusDat database, as shown in Figure 1. This figure
112 provides an insight of the range of chemicals that may be present in environmental samples.
113 To our knowledge, this is the first study "quantifying" the coverage of chemical space via
114 NTA assays.

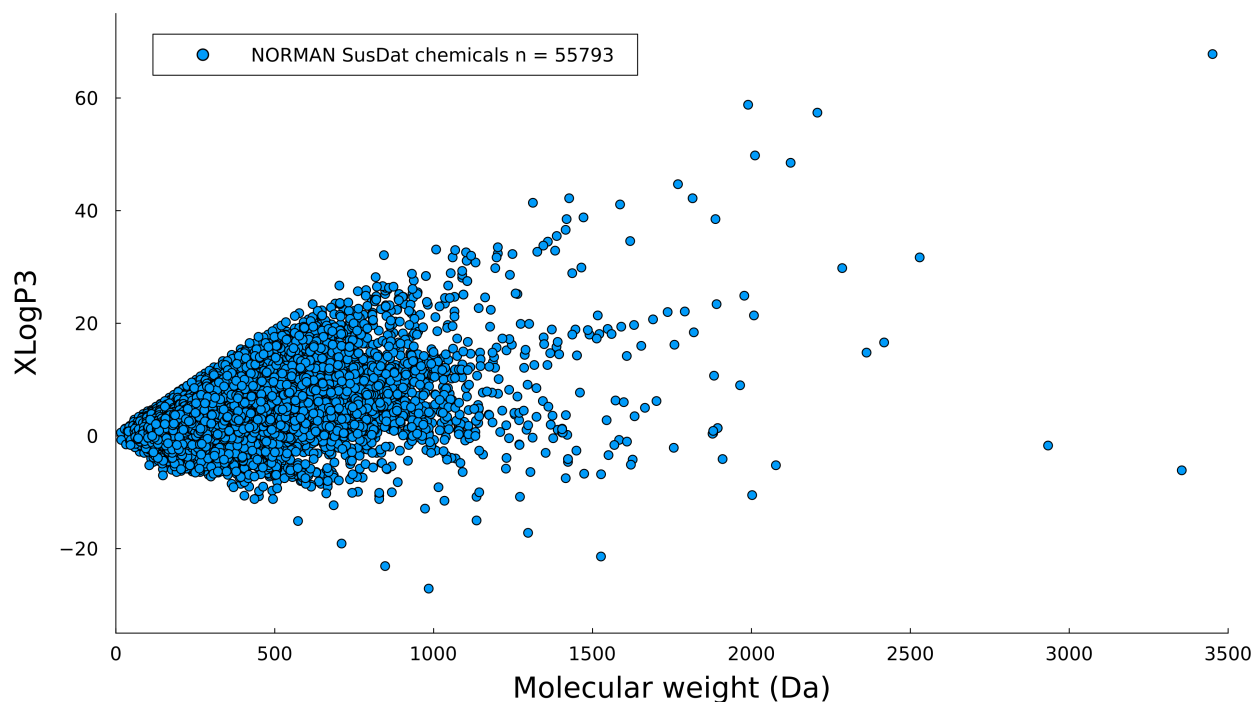


Figure 1: Distribution of all chemicals in the NORMAN SusDat database (n = 55793) based on their molecular weights (Da) and logP values.

115 Methods

116 Selection of NTA studies

117 This review is particularly focused on the development of the NTA approach in environmental
 118 studies, specifically after the discussions regarding reproducibility were initiated³⁹. Thus,
 119 we used the citation database Scopus to search for relevant studies published from 2017 to
 120 2023 in the field of non-target analysis (NTA) with a focus on environmental science. The
 121 search was limited to articles that contained the keywords "non targeted analysis", "non
 122 target analysis", "untargeted analysis", "untargeted screening", or "non-target screening"
 123 while excluding articles containing "metabolomics", "metabolic", or "gas chromatography".
 124 This initial search resulted in 377 publications adhering to the search parameters, which
 125 were then manually filtered to include only those that met a specific set of criteria.

126 The first criterion was that articles used non-target analysis to probe chemicals of emerg-
 127 ing concern, preferably in environmental matrices. Secondly, the publications had to use a

128 non-target workflow. Some articles included the desired keywords in the title or abstract but
129 were actually targeted studies with a very extensive list of target chemicals. For the same
130 reason, the third criterion was that studies used LC-HRMS for sample analysis. Studies,
131 which conducted GC-HRMS analysis were not included in the review, since such studies
132 mainly employ suspect screening rather than non-targeted analysis. Furthermore, the re-
133 cent development of NTA has been focused on LC rather than GC. Therefore, this review
134 is focused on the coverage of the chemical space by NTA conducted via LC-HRMS. Addi-
135 tionally, direct infusion studies, studies that used rare setups, or heavily modified setups
136 were excluded. Finally, review articles and studies which did not perform any identification
137 were excluded as they did not contribute any additional methods or identified compounds.
138 The search for relevant studies meeting these criteria was completed on March 1st, 2023,
139 resulting in the inclusion of 61 studies in this review⁴⁶.

140 **Collection of instrumental parameters**

141 To capture the impact of each step of the NTA workflow on chemical space coverage, we
142 extracted specific parameters used in the studies we reviewed. Sample preparation, chro-
143 matographic separation, data acquisition, and data processing were the four main steps
144 where parameters were identified. Sample preparation parameters included the sample ma-
145 trix, storage conditions, pre-storage modifications, extraction methods, and extraction con-
146 ditions where applicable. Chromatographic separation parameters included the column used,
147 eluent composition, gradient complexity, number of column volumes, column temperature,
148 and injection volume. Gradients were classified as linear, semi-linear, or complex based
149 on their complexity. The number of column volumes refers to the volume of solvent that
150 passes through a chromatography column relative to the volume of the column itself. The
151 calculation was performed using the equation 1.

$$\text{Column volumes} = \frac{F \times T \text{ run}}{\pi \times \left(\frac{dc}{2}\right)^2 \times L} \quad (1)$$

152 Where F is the flow rate (mL/min), $T \text{ run}$ is the total run time of the method (min) -
153 excluding equilibration time- dc is the internal diameter of the column (cm), and L is the
154 length of the column in (cm). HRMS instrumental parameters included the mass analyzer,
155 sampling rate (in the case of Q-TOF), resolution (in the case of Orbitrap), data acquisition
156 mode, polarity, and mass range. Data processing parameters included mass tolerance, time
157 domain alignment, mass calibration, retention time tolerance, databases used, and total
158 database size (labeled as small if ≤ 1000 compounds or large if > 1000 compounds). A
159 summary of the collected parameters can be found in Figure 2. Furthermore, we made note of
160 parameters that were not reported in order to identify which settings were commonly omitted.
161 Lastly, we gathered information on the scope of the studies. The collected parameters along
162 with the list of the publications are publicly available through this link⁴⁶.

163 **Collection of reported structures**

164 To assess the extent of chemical space coverage by recent NTA studies, we extracted the
165 reported structures. To ensure the reliability and accuracy of our analysis, we only included
166 structures identified with a high level of confidence (i.e levels one and two on the Schyman-
167 ski scale), which is less susceptible to false positive identifications⁴⁵. For each compound,
168 SMILES, IUPAC name, and the regular names provided by the authors were extracted. Fi-
169 nally, we excluded articles from our chemical space coverage assessment if the authors did
170 not specify the identification level, did not include the identified compounds in either the
171 article or supplementary materials, or only reported compounds within their target list.

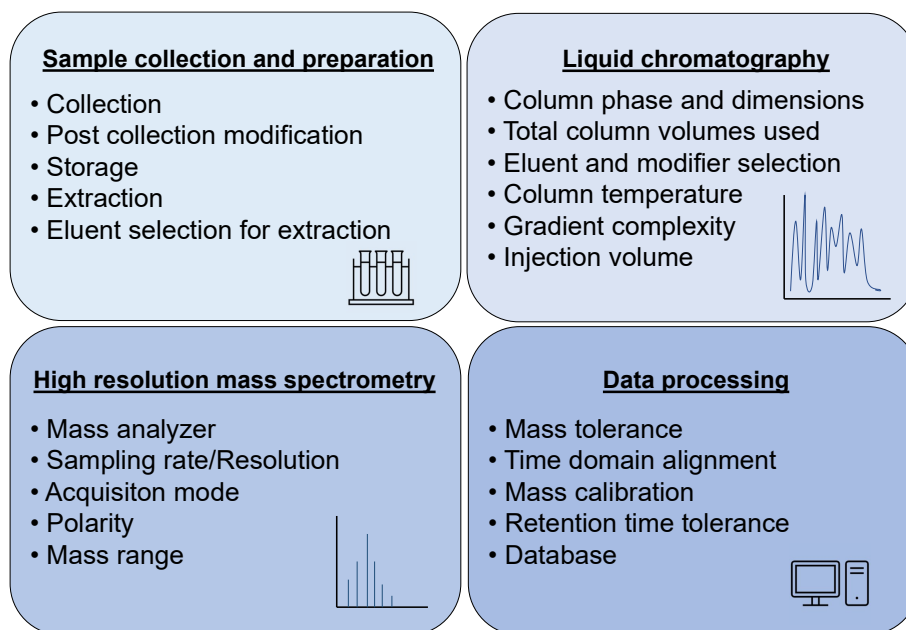


Figure 2: Summary of the main instrumental parameters collected from the reviewed NTA studies

172 Data processing

173 The list of the collected compounds was stored in CSV format, and Julia version 1.7 was
 174 used to import and process the data. A modified version of the PubChemCrawler.jl package
 175 was employed to retrieve chemical data such as XLogP3 and MW of the compounds from
 176 the PubChem database by using their available identifiers (SMILES, IUPAC, InChIKey, or
 177 regular name)⁴⁷. logP values extracted from PubChem are generated using XlogP3 with an
 178 additive model starting from a reference compound⁴⁸. Retrieved data along with the col-
 179 lected experimental parameters were combined into a dataset that included PubChem CIDs
 180 corresponding to the compounds, their logP values, molecular weights, and experimental
 181 parameters.

182 For the evaluation of the chemical space coverage, we additionally calculated elemental
 183 mass defects (EMD) of six elemental ratios (CO, CCl, CN, CS, CF, and CH) for each
 184 collected compound and the ones included in the NORMAN SusDat database⁴⁹. EMD

185 values were used to cluster structurally similar compounds together and separate others, as
186 they incorporate structural information and are used to compare compounds based on their
187 elemental composition⁵⁰. The combination of logP, MW, and EMDs was used for principal
188 component analysis (PCA), which is an unsupervised algorithm for dimensional reduction
189 combining variables into principal components⁵¹. This approach is able to identify trends and
190 clusters in the data sets. Prior to the analysis the data was mean-centered and scaled to keep
191 the initial weight of all variables comparable. PCA was performed using the ScikitLearn.jl
192 julia package and in total three principal components were utilized.

193 The NORMAN SusDat database was used for the approximation of the chemical space
194 of environmental samples. While the chemical space comprises both known and unknown
195 compounds, it is practically impossible to include the latter in our approximations. The
196 Norman SusDat database includes CECs that have either been detected in various environ-
197 mental compartments or have been identified as potential CECs, providing a comprehensive
198 set of chemicals with a wide coverage of physical and chemical properties, and structures²⁰.

199 Finally, the classes of the collected compounds were defined to illustrate the frequency
200 of identification of specific classes. To obtain the class of each CEC, the corresponding
201 InChIKey was used to generate information on superclasses, classes, and sub-classes of each
202 compound via ClassyFire. ClassyFire divides a given chemical compound into classes based
203 on its structural features (i.e. functional groups)⁵².

204 Discussion

205 In this review, we estimated the coverage of the chemical space of environmental samples
206 by investigating recent NTA studies. To evaluate the impact of selected workflow paramete-
207 ters on the coverage of chemical space, we collected information on these parameters (e.g.
208 mass analyzer, data acquisition mode, ionization mode and size of the database used) from
209 the studies⁴⁶. The identified compounds were categorized into classes and their relative

210 frequency of occurrence was determined. XLogP3, MW, and EMDs were used to represent
211 the vastness of the chemical space, approximated with the NORMAN SusDat Database.
212 PCA was employed to illustrate the coverage of the space of chemicals detected in recent
213 environmental studies.

214 **Overview of the studies**

215 In total, 61 studies were collected, with 55 of them published since 2020. Only studies
216 using NTA were included, while those using screening or targeted approaches but claiming
217 to be untargeted were excluded. This indicates that $\approx 90\%$ of the reviewed studies were
218 published in the last three years, yielding an average of more than 15 studies per year. In
219 contrast, during the period from 2017 to 2020, only six of the selected studies were published.
220 Therefore, the significant increase in the number of such studies in recent years reflects
221 the successful applications of NTA workflows in exposome analysis. The scope of these
222 studies varies, with 30 studies focusing on a wide range of chemicals and another 21 studies
223 specifically targeting groups among which are per- and polyfluoroalkyl substances (PFAS),
224 pesticides, pharmaceuticals, and illicit drugs. Such prior prioritization influences the choice
225 of experimental setup. The remaining 10 studies focused on NTA workflow development,
226 indicating a growing interest and the need for further advancements in this field.

227 **Overview of selected parameters**

228 **Sample collection and preparation**

229 The collection and preparation of samples in the non-targeted analysis (NTA) workflow can
230 introduce potential sources of loss of chemical information. Issues such as ensuring sample
231 representativeness (e.g. selecting appropriate grab or passive sampling techniques), address-
232 ing potential sample contamination, accounting for matrix effects, optimizing extraction
233 methods for selectivity, and avoiding bias towards specific chemical groups are important

234 considerations in NTA^{2,23,26}. These challenges may impact the accuracy and reliability of
235 NTA results, potentially affecting the comprehensiveness and quality of the chemical infor-
236 mation obtained from the analysis. Therefore, careful attention to sample collection and
237 preparation steps are essential to minimize potential sources of bias and ensure robust and
238 reliable NTA outcomes.

239 The majority of the collected studies (67%) analyzed water samples (n = 41). Other
240 matrices that were investigated include biota (n = 5), dust (n = 3), urine (n = 3), atmospheric
241 particulate matter (n = 2), paper (n = 2), serum (n = 2), blood (n = 1), human hair (n =
242 1), ovarian follicular fluid (n = 1), sewage sludge (n = 1), snow (n = 1), and surface soil (n
243 = 1), turtle tissue (n = 1).

244 To prevent microbiological growth, the studies on water samples reported a conservation
245 step, which involved either adding an acid or storing the sample at a temperature of -20°C or
246 4°C. Out of the 41 water studies, 7 studies either did not include a step to stop microbiological
247 growth or did not report it. If this step was omitted, it could significantly alter the sample's
248 final composition when it is eventually analyzed in the laboratory^{53,54}.

249 Around 54% of publications analyzing water included a sample filtering step prior to
250 analysis. This step is a compromise to preserve the LC system and column but may lead
251 to the loss of the chemicals adsorbed to the particle's surface. Approximately 67% of stud-
252 ies included solid phase extraction (SPE) in their sample preparation, out of which 73%
253 used reversed-phase hydrophilic-lipophilic balance (HLB) SPE. However, only 29% of stud-
254 ies with SPE used acidic and/or basic modifiers in the extraction eluents. That implies that
255 most studies using only HLB SPE are potentially leaving ionizable compounds on the sor-
256 bent and may exclude them from the analysis. The remaining studies employed alternative
257 pretreatment techniques among which vacuum-assisted evaporation, centrifugation, liquid-
258 liquid extraction (LLE), ultrasonic extraction as well as their combination. These choices are
259 mostly dictated by the sample nature/matrix. There were three studies that performed no
260 sample extraction and injected directly into the LC-MS with a higher injection volume⁵⁵⁻⁵⁷.

261 While this protocol minimizes sample adulteration and keeps the sampling of chemical space
262 more comprehensive, it can also pose a challenge to detection sensitivity due to the low
263 analyte concentration²³.

264 Overall the sample collection and preparation section is well reported in the selected
265 studies. However, many of the studies focused on analyzing a wide range of chemicals do
266 not explore alternative extraction methods to ensure a more comprehensive coverage of the
267 chemical space. This could result in a bias towards specific compounds, rather than capturing
268 a more diverse set of chemicals.

269 **Liquid chromatography**

270 Chromatographic separation is employed to minimize sample complexity by spreading ana-
271 lytes across the time axis. This helps to reduce ion suppression (matrix effect) and provides
272 additional information (retention time) for the identification of the analytes. The chemistry
273 of the stationary phase along with the elution conditions affects the quality of separation
274 and the type of analytes being retained. Thus, the selection of chromatographic conditions
275 heavily influences the coverage of the chemical space of the sample²³.

276 The majority of NTA studies use conventional reverse-phase separation with a generic
277 C18 column. Optimization of the separation includes proper selection of eluents and modi-
278 fiers, including suitable elution power and gradient setup, to avoid co-elution and excessive
279 or insufficient retention of chemicals⁵⁸. A simple linear gradient of an aqueous phase and
280 methanol or acetonitrile from low to high percentage is most widely accepted for the wide
281 scope screening. This method proved its reproducibility across different scopes of the stud-
282 ies²⁶. However, this strategy focuses on polar to semipolar compounds, potentially excluding
283 very polar (i.e. logP smaller than -2) and very hydrophobic substances (i.e. logP larger 6)
284 from the comprehensive investigation of the chemical composition of samples⁵⁹. To cover
285 the polar part of the chemical space, orthogonal methods such as hydrophilic interaction
286 chromatography (HILIC) become more popular while for hydrophobic volatile chemicals,

287 GC is a widely used technique^{55,60,61}. Finally, to ensure the reproducibility and reliability of
288 the studies parameters such as injection volume and column temperature should be properly
289 reported⁶².

290 More than 90% of the collected studies used a C18 column for the separation, among
291 which almost all were endcapped with a column length of 50mm (20%), 100mm (49%), or
292 150mm (30%). Column diameters were either 2.1mm (78% of the studies), 3mm (16%),
293 or 2mm (3%) with the particle diameter under 3.5 μm . Additionally, two different studies
294 reported 4.6mm and 0.05mm column diameters. Although applying a simple gradient ensures
295 higher reproducibility of the method, only half of the studies (replaced $\approx 51\% \approx 49\%$) used
296 a linear gradient, while around 32% used a semi-linear gradient and the remaining (18 %)
297 used a more complex type of gradient.

298 The median number of column volumes eluted in the studies is 15.9, with an interquartile
299 range of 15.9. The use of a sufficient number of column volumes should ensure the complete
300 elution of most hydrophobic compounds (high logP and MW) and the absence of carryover.
301 The optimal number depends on the stationary phase, eluent power, and analytes them-
302 selves⁶³. Nevertheless, the widely accepted hypothesis is that there is a linear relationship
303 between logP and retention/number of column volumes used. The hypothesis is applied
304 for the reverse phase mode with comparable C18 selectivity, similar gradients, and eluent
305 composition^{64,65}. However, our results do not indicate the presence of a linear relationship
306 between the number of column volumes and logP of the chemicals, since no clear linear
307 pattern could be identified between these parameters (Figure S1).

308 In addition, the column temperatures used were all slightly above room temperature
309 which is favorable for repeatability and reproducibility⁶². 31% of publications used 40°C,
310 16% used 35°C, replaced 11% used 30°C, two studies held the column at 25°C, one at 20°C,
311 one at 45°C and one at 50°C. About 29% of papers did not report the column temperature,
312 which hinders the reproducibility of experiments.

313 Finally, 18% of the studies did not report the injection volume used. Injection volume

314 should not have a large effect on the final observed chemical space as they depend on the
315 extraction method and efficiencies. Nevertheless, the success of the method's transfer de-
316 pends on it. Most of the studies used either 5 (n = 17) or 10 μL (n = 13) injection volume,
317 which is adequate when using SPE extraction. The remaining were spread across 1, 3, 4, 7,
318 20, 100, 140, and 660 μL .

319 To conclude, despite the rising discussion about reporting quality^{23, 66–68} chromatographic
320 separation parameters in the collected studies were not always properly reported. Proper
321 harmonized reporting ensures successful method transfer, whereas inconsistent reporting
322 raises questions related to the reproducibility of the study, reliability of the results, and the
323 possibility of retrospective studies. While the majority of the studies seek to comprehensively
324 investigate the chemical composition of the samples, only approximately 10% employ an
325 alternative to the conventional approach to analyze the samples. Lastly, the hypothetical
326 linear trend between logP and retention was not confirmed, indicating the need for more
327 sophisticated strategies for method development and optimization.

328 **High resolution mass spectrometry**

329 The Orbitrap and the quadrupole time of flight (QTOF) equipped with electrospray ioniza-
330 tion (ESI) are the two most commonly used HRMS instruments in liquid chromatography-
331 based (LC) NTA experiments. For complimentary analysis, it is recommended to perform
332 separate experiments in both positive and negative modes⁶⁹. The mass resolution of Orbitrap
333 mass analyzers is generally higher than that of QTOF, but both can provide high-resolution
334 mass spectra (Resolution $\geq 30,000$)⁷⁰.

335 In QTOF, resolution is determined by the architecture of the mass analyzer⁷¹, while
336 for Orbitrap, the resolving power depends on a user specified resolution. In the case of
337 Orbitrap, the speed of scans is directly related to the spectral resolution. However, the
338 increase in mass resolution is limited by the time required for scanning operations. For
339 QTOF, a crucial parameter for data quality is the sampling speed, which is reported as

340 spectra per second in Hz. If the scan rate is too high, fewer ions are sampled, which can
341 lead to a sensitivity issue. Conversely, if the scan rate is too low, fewer data points on the
342 time axis are recorded, potentially causing missed detection of analytes eluting in a narrow
343 time range⁷².

344 MS/MS spectra for structure elucidation are recorded using either data-dependent ac-
345 quisition (DDA) or data-independent acquisition (DIA). DDA mode records fragments of
346 pre-selected precursor ions (which are chosen based on their abundance or via an inclusion
347 list) while DIA mode fragments all precursor ions within a certain mass range. The latter
348 is preferable for comprehensive investigations of complex samples. However, DDA mode is
349 currently the preferred choice in environmental studies, partly due to the limited availability
350 of processing tools for DIA files and also because the DIA experimental setup is not com-
351 monly employed with Orbitrap mass analyzers²⁶. QTOF analyzers are more commonly used
352 for DIA due to higher data acquisition rates.

353 Roughly, half of the collected studies (n = 31) utilized an Orbitrap mass analyzer, while
354 the other half employed a QTOF mass analyzer. However, a significant proportion (approx-
355 imately 74%) of the studies reported using DDA, which inherently limits their results to
356 predefined ions. The scan rate for QTOF analyzers was mostly set at 4 Hz, although some
357 studies operated at lower rates of 3, 2, or 1 Hz. Many studies using Orbitrap analyzers op-
358 erated at a resolution of 70,000, while some studies used lower resolutions with a minimum
359 of 35,000 and higher with a maximum of 240,000. Approximately 22% of the studies did not
360 report either resolution or scan rate.

361 Less than half of the studies (around 42%) conducted separate experiments in positive
362 and negative modes, utilizing multiple injections, different modifiers, and sometimes different
363 columns, which is considered a more suitable scenario for achieving comprehensive coverage of
364 chemical space. In approximately 30% of the studies, MS was operated only in positive mode.
365 There were eleven publications where the analysis was reported in both modes, but the details
366 were insufficient to determine if the experiment was performed simultaneously or separately

367 in both modes. In three other studies, an exclusively negative mode was used to prioritize
368 a specific group of compounds of interest, such as PFAS⁷³⁻⁷⁵, deliberately narrowing down
369 the investigated chemical space. Finally, two of the reviewed studies employed simultaneous
370 positive and negative ionization modes with formic acid as a modifier. This approach is
371 not preferable for NTA given that acidic additives are not always the optimal for a negative
372 ionization mode. Additionally, the acquired data becomes extremely complex and lacks
373 quality for reliable and robust processing.

374 The selected mass range in the collected studies is between 50-1200 m/z, which is based
375 on approximated chemical space covering the largest part. However, some studies set their
376 maximum m/z at 1000 or lower, which leads to the exclusion of the part of chemical space
377 with higher MW.

378 To conclude, despite recent advancements in DIA technology, DDA remains the predom-
379 inant choice in the reviewed studies. However, the recommended approach for improved
380 reproducibility and reliability of NTA studies, and to enhance coverage of chemical space
381 in environmental and metabolomics research, is to acquire data in DIA mode for initial
382 screening and then continue with DDA for individual feature identification^{23,76-78}. Finally,
383 in terms of reproducibility the lack of comprehensively reported information hinders method
384 transfer and therefore it warrants actions towards a harmonized reporting strategy⁶⁶⁻⁶⁸.

385 **Data processing**

386 Data processing is considered a major bottleneck in NTA workflows. It refers to a series of
387 procedures that starts with the data conversion and ends with the feature identification²³.
388 One of the steps for reliable processing is the mass calibration, either external or internal.
389 During this step the measured m/z values of known structures are compared against theoret-
390 ical m/z values. These shifts/correction factors are applied to all mass channels, depending
391 on the instrumental setup. This step ensures the quality of the spectra in terms of accu-
392 rate mass measurement⁷⁹. An inadequate mass calibration may result in false positive and

393 negative detections during the identification⁸⁰.

394 One of the last steps of CEC identification is the use of a database to relate the MS
395 output to a known chemical structure. To proceed with the identification, experimental data
396 undergoes pre-processing steps: data compression, to remove noise and blank peaks, feature
397 detection, to find features in 3-dimensional data, componentization, to group fragments and
398 isotopologues belonging to the same compound, and feature prioritization to reduce the
399 number of irrelevant features⁸¹. Since most of the collected studies used vendor software
400 for the latter four steps, which makes it almost impossible to retrieve the information of
401 algorithms utilized, these parameters cannot be adequately discussed for their influence on
402 the coverage of chemical space. For the identification of known unknowns, pre-processed
403 data is compared with chemical databases and matched against references from available
404 spectral libraries, utilizing a combination of features, retention time, accurate mass, and
405 fragmentation pattern⁴⁵. The mass tolerance is the initial parameter used for the candidates'
406 list compilation. This parameter, along with the database used, heavily affects the results of
407 the candidate search. The number of chemicals included in databases used in the evaluated
408 studies differs from a few hundred structures in in-house libraries⁸² to tens and hundreds of
409 thousands in publicly available libraries²⁷ such as NORMAN,²⁰ MassBank⁸³ or PubChem¹⁹.
410 These search algorithms result in a set of candidate structures that ultimately must be
411 confirmed via either reference standard and/or an orthogonal method⁴⁵. The retention
412 tolerance is applied for level 1 confirmation employing either predicted or measured retention
413 times.²³

414 For the transparency and reproducibility of the method, proper reporting of applied se-
415 tups for each data processing step is essential. Nevertheless, a significant part of the studies
416 did not provide sufficient information to reproduce the results. Specifically, approximately
417 39% did not mention anything about mass calibration, while 25% reported that they per-
418 formed calibration but did not describe the procedure. Only about 36% included a report
419 on the mass calibration procedure. A large number of the papers (43%) also did not report

420 whether a retention alignment was performed. 34% did report the fact that a retention
421 alignment was done but did not specify the algorithm that was used or provided the details
422 on the parameters used. The remaining 23% of publications did report both the fact that
423 one was performed and which algorithm was used.

424 In contrast, mass tolerance applied for the search was reported in almost all studies,
425 around 95%. Among which around 76% used a mass tolerance for the database query of
426 5ppm, which is highly common in the NTA database search workflows. There were also
427 studies that used a relatively high mass tolerance of 20 ppm, 17 ppm, or 10 ppm and some
428 studies that used mass tolerances lower than 5 at 3ppm, 2ppm, and even 1ppm. Generally,
429 the studies that were using the lower mass tolerances for the database search reported a
430 higher resolution of the mass analyzer. However, there was no clear indication of whether
431 the mass tolerance applied to formula assignment or structural identification. On the other
432 hand, retention tolerances had much lower reporting rates as 45% of the studies did not
433 include this information. The remaining studies used tolerances in a range between 0.1 min
434 and 0.5 min. However, there are a few publications that used a wider tolerance, up to 1.8
435 min, which may result in a high false positive rate. Finally, approximately 9% did not report
436 the databases used or referred to the software but not the databases that the software was
437 using. The majority, 82%, used a total database size containing more than 5000 compounds,
438 while only 5 studies used databases with less than one thousand compounds.

439 The data processing step is one of the main bottlenecks for the NTA approach and thus
440 requires greater attention within the community. Nevertheless, the reporting quality needs
441 improvement. Furthermore, it was found that around 70% of the identified chemicals are
442 available in MassBank EU. That means that roughly 30% of the HRMS spectra acquired for
443 the identified compounds have not been deposited in public databases such as MassBank.
444 For NTA to reach its full potential, the expansion of publicly available spectral databases is
445 vital for the improvement of the coverage of chemical space at the identification step.

446 Explored chemical space

447 The studies yielded a total of 2657 compounds reported in the identification level 1 up to
448 2b. The contribution to the total number from each study varies between 1 and 370, with
449 a mean and median of 50 and 30, respectively (Figure S2). Among these, 1606 compounds
450 were identified as unique structures, accounting for $\approx 60\%$ of the total number of retrieved
451 compounds. This finding implies that around half of the overall variety of compounds were
452 detected more than once in various environmental compartments. However, in 7 studies,
453 there was no report of either identification level, or any identifiers, which hinders the retrieval
454 of the compounds from these studies. The class of each collected CEC was obtained and
455 displayed in Figure 3. The most commonly found compounds were benzoids, followed by
456 organocyclic compounds and then organic acids and derivatives. The latter category, along
457 with organohalogen compounds, constitutes PFAS, which have been of particular interest
458 in recent years. The median molecular weights of compounds from SusDat were 239 Da
459 and 261 Da for those collected from the studies, with a median XLogP3 of 3.2 for SusDat
460 and 2.2 for collected compounds. Based on histograms in Figure 4, compounds with the
461 most frequently occurring properties are being identified in recent NTA studies, which can
462 be partially explained by the generalized experimental workflows with reverse phase C18
463 columns.

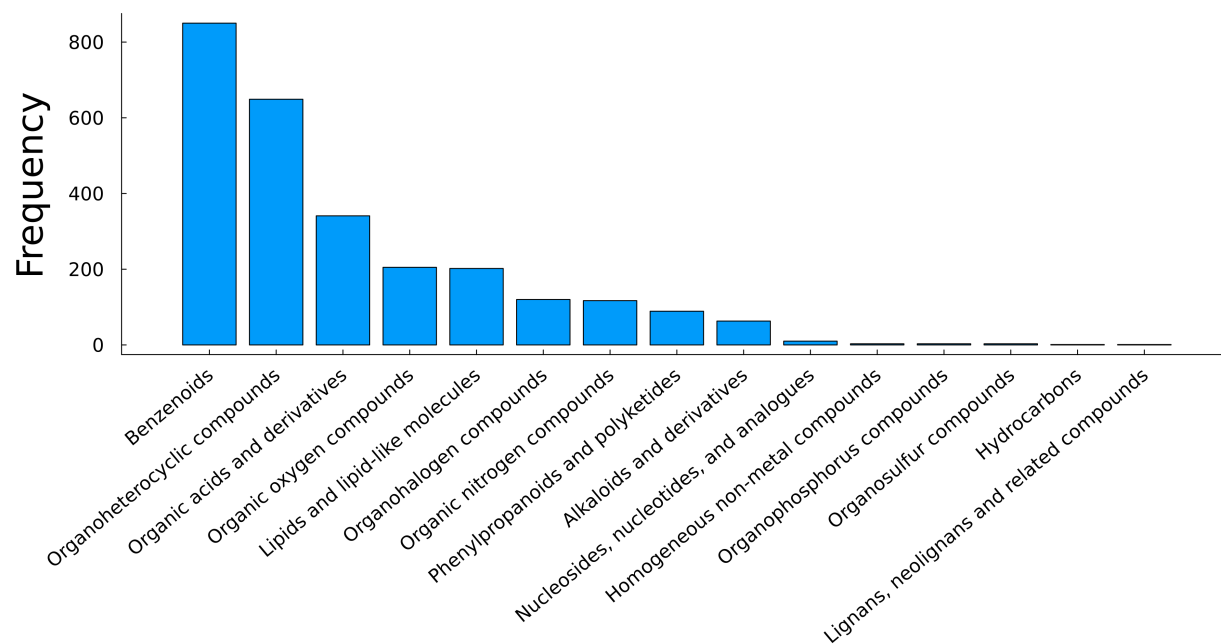


Figure 3: Histogram of all of the classes obtained from the Classyfire search for the detected CECs in reviewed studies

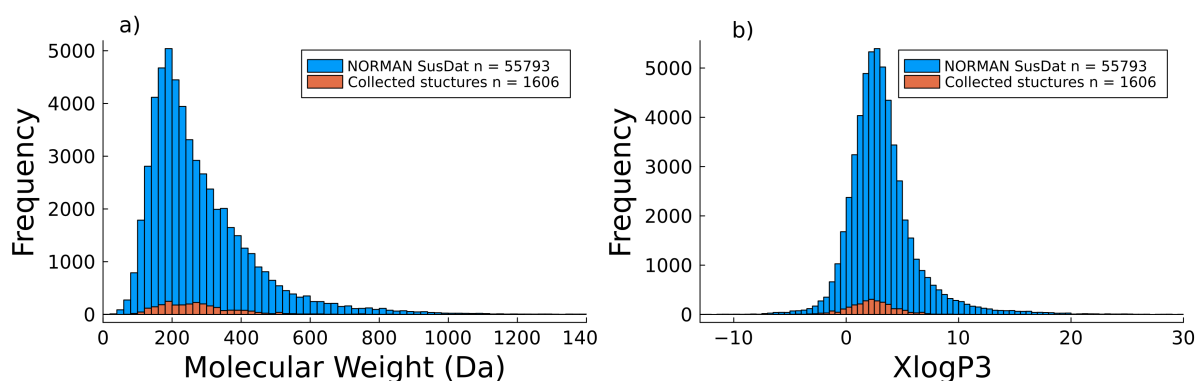


Figure 4: Molecular weights (a) and logP (b) distributions for the collected compounds (orange) and ones included in NORMAN SusDat database (blue).

464 Most of the compounds detected in the studies clustered closely together, with only a few
 465 compounds found further away from this main cluster (Figure 5). The collected compounds
 466 were analyzed in relation to their properties and plotted on a chemical space approximation
 467 represented by the NORMAN SusDat database. Figure 5 shows the plot in dimensions of

468 molecular weight (MW) and XLogP3, which emphasizes the limited space that is currently
469 explored using current non-target analysis workflows. To examine the effect of some of the
470 mass spectrometry (MS) parameters used on the explored chemical space, all compounds
471 were plotted and clustered based on factors such as the mass analyzer used, acquisition
472 mode, ionization mode, and the total database size used (Figure S3-S6). However, neither of
473 these parameters showed an unambiguous influence on the coverage of the chemical space.
474 It should be noted that the representation in MW and logP dimensions does not provide
475 information about the elemental composition of compounds or their classes, which may result
476 in an over-representation of the covered chemical space. Therefore, it is important to consider
477 other parameters beyond MW and logP when evaluating the coverage of the chemical space
478 by the collected structures.

479 The PCA scores plot in Figure 6 reveals that many regions of the chemical space are
480 unexplored. The PCA was applied to the dataset combining the collected compounds with
481 the ones from the NORMAN SusDat, with MW, XLogP3, and the EMDs as input variables.
482 The first two principal components in the analysis were found to be primarily influenced
483 by the elemental mass differences (EMDs) associated with compounds containing chlorine
484 (Cl), fluorine (F), cyanide (CN), and sulfur (S). These EMDs represent the high variability
485 in the elemental composition of the compounds and were identified as the most important
486 variables in the PCA. This indicates that fewer compounds in the dataset contain halogens,
487 nitrogen (N), and sulfur, while hydrogen (H), which is present in nearly every compound,
488 does not contribute significantly to the variability in the data. The third principal com-
489 ponent is primarily influenced by MW and XlogP3 (Figure 76). In total, the first three
490 principal components explain 74% of the variance (Figure S8). In Figure S9-S11, the cov-
491 erage of chemical space by different compound classes is displayed. Figure S11 specifically
492 highlights the coverage by organic acids and derivatives as well as organohalogen compounds.
493 The majority of PFAS, not exclusively, fall into these classes. The figures reveal that the
494 distribution of compound classes across the chemical space is not homogeneous, suggesting

495 an over-representation of certain classes of compounds. This observation can be attributed
496 to the prior prioritization of specific classes, which may bias the identification towards those
497 classes of compounds.

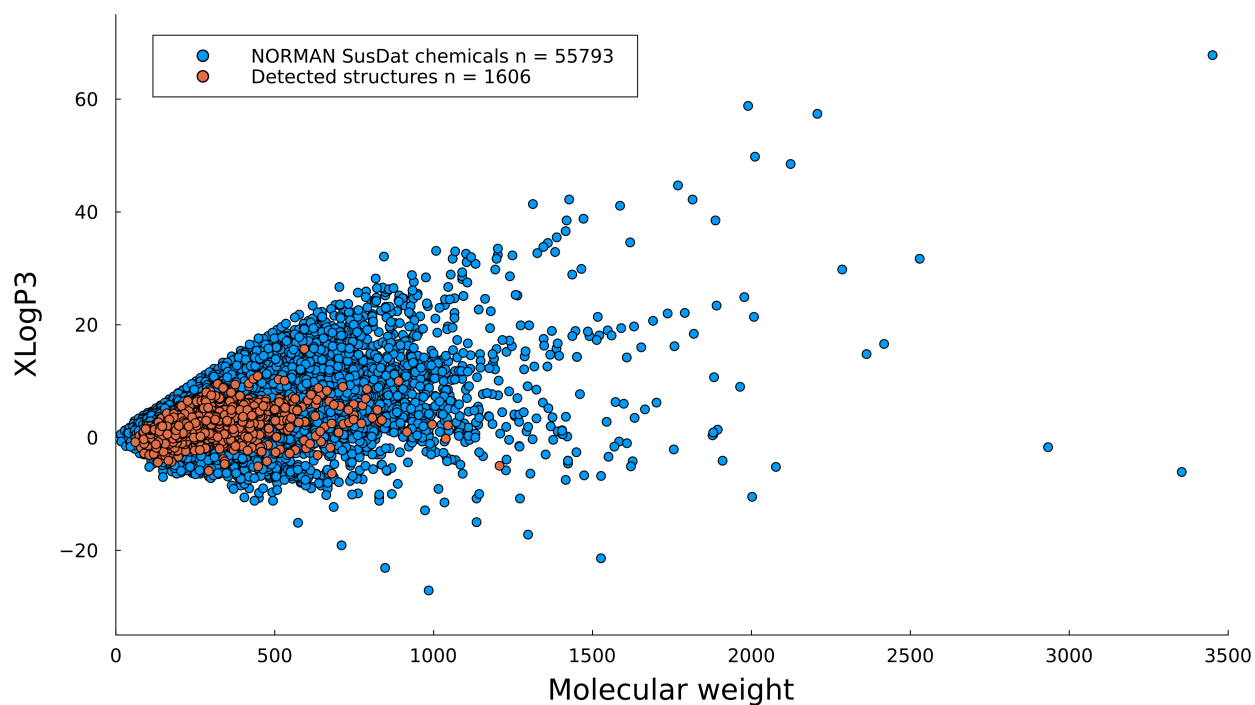


Figure 5: Distribution of all chemicals found in the reviewed articles at level 1 to 2b (orange) overlaid on NORMAN SusDat database chemicals (blue) based on their molecular weights and XlogP3 value

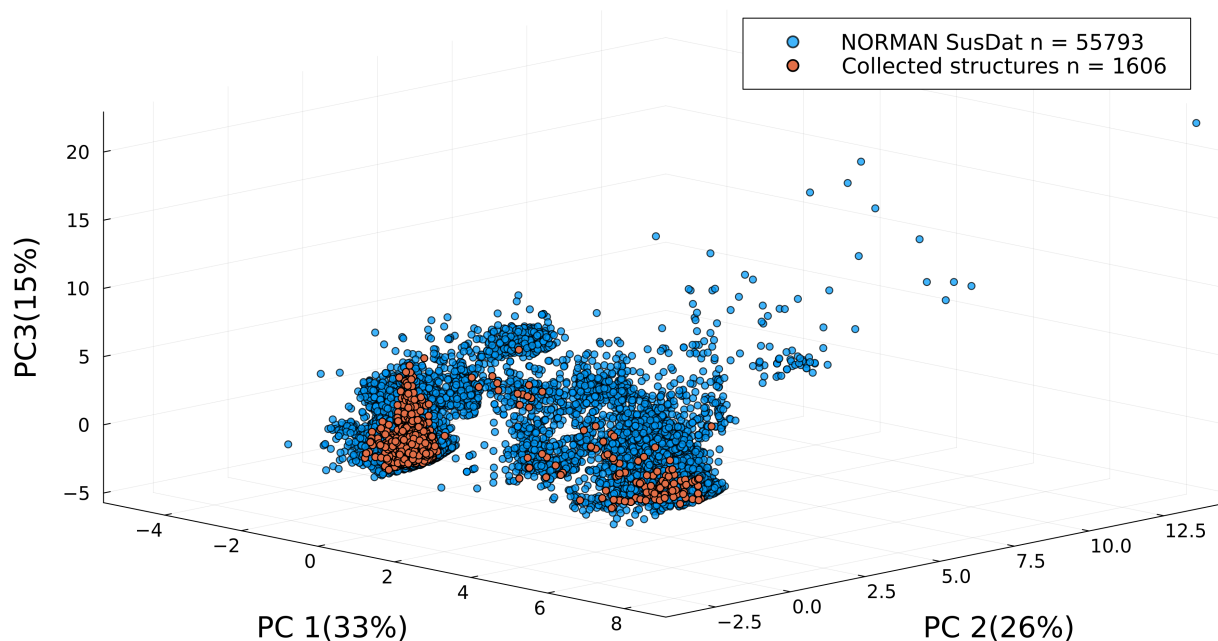


Figure 6: Scores plot of three principal components of the NORMAN SusDat database (blue) and the collected structures (orange).

498 Overall only around 2% of the estimated chemical space was covered by NTA studies
 499 investigated in this review. The coverage was defined as the number of unique structures
 500 retrieved from the reviewed studies versus the number of structures in NORMAN SusDat
 501 database as an approximation of chemical space. We used NORMAN SusDat as an approx-
 502 imation of chemical space as it contains a set of highly relevant and curated structures for
 503 environmental and human exposome. It should be noted that this is a small subspace of
 504 the total chemical space and serves as means of approximation for the true chemical space
 505 of exposome. No clear relationship between experimental conditions and coverage of the
 506 chemical space was discovered, which may indicate that the used experimental approaches
 507 are generic enough for the NTA assays. On the other hand, this may be caused by the lack
 508 of detailed and standardized reporting of the experimental conditions. Therefore, a more
 509 rigorous investigation of the parameters and standardization of reporting criteria has to be
 510 designed and performed. Although the most widely accepted properties of compounds such
 511 as logP and MW are widely used while discussing chemical space²³, in this study we showed

512 that they may not be the most relevant markers for assessing the coverage of chemical space.
513 Finally, such a low coverage emphasizes the need for more comprehensive approaches to ex-
514 perimental and data processing workflows in order to explore a broader range of the chemical
515 space and ultimately protect human and environmental health.

516 **Recommendations and Outlook**

517 Despite the ability of NTA to provide holistic information about the chemical composition
518 of the samples, their true coverage of the chemical space has not been investigated. Further-
519 more, the NTA studies have suffered from issues related to their reproducibility, due to the
520 complexity of both experimental and computational approaches employed in NTA assays.
521 One of the main bottlenecks for a more reproducible NTA assay is the lack of standardization
522 of the reporting criteria (including the experimental conditions). Our detailed investigation
523 of the previously published NTA studies further suggests the need for such criteria. Mini-
524 mum accepted experimental criteria and data processing parameters should be reported to
525 ensure the transparency and reliability of the results. The utilization of harmonized report-
526 ing tools such as BP4NTA SRT or NORMAN suspect screening reporting tools can help the
527 reproducibility and transparency of future NTA studies^{68,77,78}. This will potentially lead to
528 the acceptance of the NTA approach by the regulatory bodies.

529 The potential coverage of the chemical space should be assessed during the design of the
530 experimental setups. Most of the recent studies focused their experimental setups based
531 on the conventional workflow including HLB SPE for sample preparation, reverse phase
532 separation with C18 columns, and DDA acquisition mode, without considering alternative
533 approaches. The best practice would be an application of alternative extraction methods,
534 implementation of orthogonal techniques (e.g. RPLC and HILIC), DIA acquisition mode
535 as the first screening approach, and the application of reliable/robust data processing tools,
536 preferably open source/access. For the identification part of the workflow, the sharing of

537 experimental mass spectra of identified compounds along with their acquisition conditions
538 is vital to the progress of the community. Additionally, archiving the raw data in public
539 repositories for both the retrospective analysis as well as data processing tool development
540 is highly essential.

541 To our knowledge, no other study has evaluated the coverage of the chemical space via
542 NTA studies in such detail. However, due to the lack of standardized reporting criteria,
543 the direct impact of different experimental choices on the covered chemical space could not
544 be established. Also, our study is limited to the works published after 2017 and we only
545 included studies with clear level 1 and 2 identification reporting. Furthermore, we excluded
546 the suspect screening studies, which may result in an underestimation of the coverage of
547 NTA studies. However, our study, even though limited, clearly shows the shortcomings of
548 the current NTA practices and the need for further development in different areas - including
549 experimental setup.

550 **Acknowledgement**

551 The authors are thankful to the members of Environmental Modeling & Computational
552 Mass Spectrometry (www.emcms.info). S.S. and V.T. are thankful to the UvA Data Science
553 Center and ChemistryNL TKI for the financial support (projects Edified and SCOPE). The
554 authors thank Denice van Herwerden for her help in setting up the calculations of elemental
555 mass defects. J.W.O. is the recipient of a National Health and Medical Research Council
556 (NHMRC) Investigator Grant (EL12009209) funded by the Australian Government. K.V.T,
557 and S.S. also acknowledge funding support from the Australian National Health and Medical
558 Research Council (NHMRC; APP1185347). The Queensland Alliance for Environmental
559 Health Sciences (QAEHS), The University of Queensland (UQ), gratefully acknowledges the
560 financial support of Queensland Health. This research was funded in whole or in part by
561 NHMRC Emerging Leadership Fellowship EL1 2009209. For the purposes of open access,

562 the author has applied a Creative Commons Attribution (CC BY) license to any Author
563 Accepted Manuscript version arising from this submission.

564 Notes

565 Information retrieved in this study can be found at <https://doi.org/10.5281/zenodo>
566 [.7774345](https://doi.org/10.5281/zenodo.7774345). References to the reviewed studies and collected experimental parameters are
567 at All experimental parameters.xlsx. The script to perform the calculations is available at
568 <https://github.com//tobihul//Code-for-Critical-assessment-of-covered-chemical-space-with-LC-HRMS-non-targeted-analysis>. PubChemCrawler package is
569 available at <https://github.com/JuliaHealth/PubChemCrawler.jl>.
570

571 Supporting Information Available

572 The Supporting Information with figures (S1 - S10) showing the relationship between exper-
573 imental parameters and the covered chemical space is available at XXX.

574 References

- 575 (1) Wild, C. P. Complementing the Genome with an “Exposome”: The Outstanding Chal-
576 lenge of Environmental Exposure Measurement in Molecular Epidemiology. *Cancer*
577 *Epidemiology, Biomarkers & Prevention* **2005**, *14*, 1847–1850.
- 578 (2) Milman, B. L.; Zhurkovich, I. K. The chemical space for non-target analysis. *Trends*
579 *Analyt. Chem.* **2017**, *97*, 179–187.
- 580 (3) Wambaugh, J. F. et al. New approach methodologies for exposure science. *Curr. Opin.*
581 *Toxicol.* **2019**, *15*, 76–92.

- 582 (4) Escher, B. I.; Stapleton, H. M.; Schymanski, E. L. Tracking complex mixtures of chem-
583 icals in our changing environment. *Science* **2020**, *367*, 388–392.
- 584 (5) Reymond, J. L. The Chemical Space Project. *Acc. Chem. Res.* **2015**, *48*, 722–730.
- 585 (6) Van Deursen, R.; Reymond, J. L. Chemical space travel. *ChemMedChem* **2007**, *2*,
586 636–640.
- 587 (7) Reymond, J. L.; Awale, M. Exploring chemical space for drug discovery using the
588 chemical universe database. *ACS Chem. Neurosci.* **2012**, *3*, 649–657.
- 589 (8) McEachran, A. D.; Sobus, J. R.; Williams, A. J. Identifying known unknowns using the
590 US EPA’s CompTox Chemistry Dashboard. *Anal. Bioanal. Chem.* **2017**, *409*, 1729–
591 1735.
- 592 (9) Vermeulen, R.; Schymanski, E. L.; Barabási, A.-L.; Miller, G. W. The exposome and
593 health: Where chemistry meets biology. *Science* **2020**, *367*, 392–396.
- 594 (10) Landrigan, P. J. et al. The Lancet Commission on pollution and health. *Lancet* **2018**,
595 *391*, 462–512.
- 596 (11) Petrie, B.; Barden, R.; Kasprzyk-Hordern, B. A review on emerging contaminants in
597 wastewaters and the environment: Current knowledge, understudied areas and recom-
598 mendations for future monitoring. *Water Research* **2015**, *72*, 3–27.
- 599 (12) Samanipour, S.; O’Brien, J. W.; Reid, M. J.; Thomas, K. V.; Praetorius, A. From
600 Molecular Descriptors to Intrinsic Fish Toxicity of Chemicals: An Alternative Approach
601 to Chemical Prioritization. *Environ. Sci. Technol.* **2022**,
- 602 (13) Dulio, V.; van Bavel, B.; Brorström-Lundén, E.; Harmsen, J.; Hollender, J.;
603 Schlabach, M.; Slobodnik, J.; Thomas, K.; Koschorreck, J. Emerging pollutants in
604 the EU: 10 years of NORMAN in support of environmental policies and regulations.
605 *Environ Sci Eur* **2018**, *30*.

- 606 (14) Williams, A. J.; Grulke, C. M.; Edwards, J.; McEachran, A. D.; Mansouri, K.;
607 Baker, N. C.; Patlewicz, G.; Shah, I.; Wambaugh, J. F.; Judson, R. S.; Richard, A. M.
608 The CompTox Chemistry Dashboard: A community data resource for environmental
609 chemistry. *J Cheminform* **2017**, *9*.
- 610 (15) Richard, A. M. et al. ToxCast Chemical Landscape: Paving the Road to 21st Century
611 Toxicology. *Chem. Res. Toxicol.* **2016**, *29*, 1225–1251.
- 612 (16) Tian, Z. et al. A ubiquitous tire rubber–derived chemical induces acute mortality in
613 coho salmon. *Science* **2021**, *371*, 185–189.
- 614 (17) Sauv e, S.; Desrosiers, M. A review of what is an emerging contaminant. *Chem. Cent.*
615 *J.* **2014**, *8*, 15.
- 616 (18) Maddela, N. R.; Ramakrishnan, B.; Kakarla, D.; Venkateswarlu, K.; Megharaj, M.
617 Major contaminants of emerging concern in soils: a perspective on potential health
618 risks. *RSC Advances* **2022**, *12*, 12396–12415.
- 619 (19) Wang, Y.; Xiao, J.; Suzek, T. O.; Zhang, J.; Wang, J.; Bryant, S. H. PubChem: a
620 public information system for analyzing bioactivities of small molecules. *Nucleic Acids*
621 *Res.* **2009**, *37*, 623–633.
- 622 (20) Mohammed Taha, H. et al. The NORMAN Suspect List Exchange (NORMAN-SLE):
623 facilitating European and worldwide collaboration on suspect screening in high resolu-
624 tion mass spectrometry. *Environ Sci Eur* **2022**, *34*.
- 625 (21) Ace a, J.; Stampachiachiere, S.; P rez, S.; Barcel , D. Advances in liquid chromato-
626 graphy - High-resolution mass spectrometry for quantitative and qualitative environmental
627 analysis. *Anal. Bioanal. Chem.* **2015**, *407*, 6289–6299.
- 628 (22) Pic , Y.; Barcel , D. Transformation products of emerging contaminants in the envi-

- 629 ronment and high-resolution mass spectrometry: A new horizon. *Anal Bioanal Chem*
630 **2015**, *407*, 6257–6273.
- 631 (23) Schulze, B.; Jeon, Y.; Kaserzon, S.; Heffernan, A. L.; Dewapriya, P.; O'Brien, J.;
632 Gomez Ramos, M. J.; Ghorbani Gorji, S.; Mueller, J. F.; Thomas, K. V.; Samanipour, S.
633 An assessment of quality assurance/quality control efforts in high resolution mass spec-
634 trometry non-target workflows for analysis of environmental samples. *TrAC Trends in*
635 *Analytical Chemistry* **2020**, *133*, 116063.
- 636 (24) Muir, D. C.; Howard, P. H. Are there other persistent organic pollutants? A challenge
637 for environmental chemists. *Environ. Sci. Technol.* **2006**, *40*, 7157–7166.
- 638 (25) Samanipour, S.; Martin, J. W.; Lamoree, M. H.; Reid, M. J.; Thomas, K. V. Letter to
639 the Editor: Optimism for Nontarget Analysis in Environmental Chemistry. *Environ.*
640 *Sci. Technol.* **2019**, *53*, 5529–5530.
- 641 (26) Menger, F.; Gago-Ferrero, P.; Wiberg, K.; Ahrens, L. Wide-scope screening of polar
642 contaminants of concern in water: A critical review of liquid chromatography-high
643 resolution mass spectrometry-based strategies. *Trends Environ. Anal. Chem.* **2020**, *28*,
644 e00102.
- 645 (27) Hollender, J.; Schymanski, E. L.; Singer, H. P.; Ferguson, P. L. Nontarget Screening
646 with High Resolution Mass Spectrometry in the Environment: Ready to Go? *Environ.*
647 *Sci. Technol.* **2017**, *51*, 11505–11512.
- 648 (28) Zedda, M.; Zwiener, C. Is nontarget screening of emerging contaminants by LC-HRMS
649 successful? A plea for compound libraries and computer tools. *Anal. Bioanal. Chem.*
650 **2012**, *403*, 2493–2502.
- 651 (29) Krauss, M.; Singer, H.; Hollender, J. LC-high resolution MS in environmental analysis:
652 From target screening to the identification of unknowns. *Anal. Bioanal. Chem.* **2010**,
653 *397*, 943–951.

- 654 (30) McCord, J. P.; Groff, L. C.; Sobus, J. R. Quantitative non-targeted analysis: Bridging
655 the gap between contaminant discovery and risk characterization. *Environ. Int.* **2022**,
656 *158*, 107011.
- 657 (31) Schmidt, T. C. Recent trends in water analysis triggering future monitoring of organic
658 micropollutants. *Anal. Bioanal. Chem.* **2018**, *410*, 3933–3941.
- 659 (32) Samanipour, S.; Reid, M. J.; Thomas, K. V. Statistical Variable Selection: An Alterna-
660 tive Prioritization Strategy during the Nontarget Analysis of LC-HR-MS Data. *Anal.*
661 *Chem.* **2017**, *89*, 5585–5591.
- 662 (33) van Herwerden, D.; O'Brien, J. W.; Lege, S.; Pirok, B. W. J.; Thomas, K. V.; Sama-
663 nipour, S. Cumulative Neutral Loss Model for Fragment Deconvolution in Electrospray
664 Ionization High-Resolution Mass Spectrometry Data. *Anal. Chem.* **2023**,
- 665 (34) Hohrenk, L. L.; Itzel, F.; Baetz, N.; Tuerk, J.; Vosough, M.; Schmidt, T. C. Comparison
666 of Software Tools for Liquid Chromatography-High-Resolution Mass Spectrometry Data
667 Processing in Nontarget Screening of Environmental Samples. *Anal. Chem.* **2020**, *92*,
668 1898–1907.
- 669 (35) Samanipour, S.; Baz-Lomba, J. A.; Alygizakis, N. A.; Reid, M. J.; Thomaidis, N. S.;
670 Thomas, K. V. Two stage algorithm vs commonly used approaches for the suspect
671 screening of complex environmental samples analyzed via liquid chromatography high
672 resolution time of flight mass spectroscopy: A test study. *J. Chromatograph A* **2017**,
673 *1501*, 68–78.
- 674 (36) Hernández, F. et al. The role of analytical chemistry in exposure science: Focus on the
675 aquatic environment. *Chemosphere* **2019**, *222*, 564–583.
- 676 (37) Brack, W.; Hollender, J.; de Alda, M. L.; Müller, C.; Schulze, T.; Schymanski, E.;
677 Slobodnik, J.; Krauss, M. High-resolution mass spectrometry to complement monitor-

- 678 ing and track emerging chemicals and pollution trends in European water resources.
679 *Environ Sci Eur* **2019**, *31*, 62.
- 680 (38) Black, G. et al. Exploring chemical space in non-targeted analysis: a proposed
681 ChemSpace tool. *Anal. Bioanal. Chem.* **2023**, *415*, 35–44.
- 682 (39) Hites, R. A.; Jobst, K. J. Is Nontargeted Screening Reproducible? *Environ. Sci. Tech-*
683 *nol.* **2018**, *52*, 11975–11976.
- 684 (40) Ulrich, E. M.; Sobus, J. R.; Grulke, C. M.; Richard, A. M.; Newton, S. R.; Stry-
685 nar, M. J.; Mansouri, K.; Williams, A. J. EPA’s non-targeted analysis collaborative
686 trial (ENTACT): genesis, design, and initial findings. *Anal. Bioanal. Chem.* **2019**, *411*,
687 853–866.
- 688 (41) Fisher, C. M.; Peter, K. T.; Newton, S. R.; Schaub, A. J.; Sobus, J. R. Approaches for
689 assessing performance of high-resolution mass spectrometry-based non-targeted anal-
690 ysis methods. *Anal. Bioanal. Chem.* **2022**, *414*, 6455–6471.
- 691 (42) Knolhoff, A. M.; Premo, J. H.; Fisher, C. M. A proposed quality control standard
692 mixture and its uses for evaluating nontargeted and suspect screening LC/HR-MS
693 method performance. *Anal. Chem.* **2021**, *93*, 1596–1603.
- 694 (43) Singh, R. R.; Chao, A.; Phillips, K. A.; Xia, X. R.; Shea, D.; Sobus, J. R.; Schy-
695 manski, E. L.; Ulrich, E. M. Expanded coverage of non-targeted LC-HRMS using at-
696 mospheric pressure chemical ionization: a case study with ENTACT mixtures. *Anal.*
697 *Bioanal. Chem.* **2020**, *412*, 4931–4939.
- 698 (44) Newton, S. R.; Sobus, J. R.; Ulrich, E. M.; Singh, R. R.; Chao, A.; McCord, J.;
699 Laughlin-Toth, S.; Strynar, M. Examining NTA performance and potential using for-
700 tified and reference house dust as part of EPA’s Non-Targeted Analysis Collaborative
701 Trial (ENTACT). *Anal. Bioanal. Chem.* **2020**, *412*, 4221–4233.

- 702 (45) Schymanski, E. L.; Jeon, J.; Gulde, R.; Fenner, K.; Ruff, M.; Singer, H. P.; Hollender, J.
703 Identifying small molecules via high resolution mass spectrometry: Communicating
704 confidence. *Environ. Sci. Technol.* **2014**, *48*, 2097–2098.
- 705 (46) Hulleman, T.; Turkina, V.; O'Brien, J.; Chojnacka, A.; Thomas, K. V.; Saminopour, S.
706 Data files for: Critical assessment of covered chemical space with LC-HRMS non-
707 targeted analysis. 2022; <https://doi.org/10.5281/zenodo.7774345>.
- 708 (47) Hulleman, T. Code for: Critical assessment of covered chemical space with LC-HRMS
709 non-targeted analysis. 2023; [https://github.com/tobihul/Code-for-Critical-a-
710 sssessment-of-covered-chemical-space-with-LC-HRMS-non-targeted-analysis](https://github.com/tobihul/Code-for-Critical-assessment-of-covered-chemical-space-with-LC-HRMS-non-targeted-analysis).
- 711 (48) Cheng, T.; Zhao, Y.; Li, X.; Lin, F.; Xu, Y.; Zhang, X.; Li, Y.; Wang, R.; Lai, L.
712 Computation of octanol-water partition coefficients by guiding an additive model with
713 knowledge. *J. Chem. Inf. Model.* **2007**, *47*, 2140–2148.
- 714 (49) van Herwerden, D.; O'Brien, J. W.; Choi, P. M.; Thomas, K. V.; Schoenmakers, P. J.;
715 Samanipour, S. Naive Bayes classification model for isotopologue detection in LC-
716 HRMS data. *Chemom. Intell. Lab. Syst.* **2022**, *223*, 104515.
- 717 (50) Sleno, L. The use of mass defect in modern mass spectrometry. *J. Mass. Spectrom.*
718 **2012**, *47*, 226–236.
- 719 (51) Hastie, T.; Tibshirani, R.; Friedman, J. *The Elements of Statistical Learning Data
720 Mining, Inference, and Prediction*; Springer, 2009.
- 721 (52) Djoumbou Feunang, Y.; Eisner, R.; Knox, C.; Chepelev, L.; Hastings, J.; Owen, G.;
722 Fahy, E.; Steinbeck, C.; Subramanian, S.; Bolton, E.; Greiner, R.; Wishart, D. S. Classy-
723 Fire: automated chemical classification with a comprehensive, computable taxonomy.
724 *J. Cheminform.* **2016**, *8*, 61.

- 725 (53) Krueve, A. Semi-quantitative non-target analysis of water with liquid
726 chromatography/high-resolution mass spectrometry: How far are we? *Rapid*
727 *Commun. Mass Spectrom.* **2019**, *33*, 54–63.
- 728 (54) Lyytikäinen, M.; Kukkonen, J. V.; Lydy, M. J. Analysis of pesticides in water and
729 sediment under different storage conditions using gas chromatography. *Arch. Environ.*
730 *Contam. Toxicol.* **2003**, *44*, 437–444.
- 731 (55) Been, F.; Krueve, A.; Vughs, D.; Meekel, N.; Reus, A.; Zwartsen, A.; Wessel, A.; Fis-
732 cher, A.; ter Laak, T.; Brunner, A. M. Risk-based prioritization of suspects detected
733 in riverine water using complementary chromatographic techniques. *Water Res.* **2021**,
734 *204*, 117612.
- 735 (56) Hu, L. X.; Olaitan, O. J.; Li, Z.; Yang, Y. Y.; Chimezie, A.; Adepoju-Bello, A. A.;
736 Ying, G. G.; Chen, C. E. What is in Nigerian waters? Target and non-target screening
737 analysis for organic chemicals. *Chemosphere* **2021**, *284*, 131546.
- 738 (57) Köppe, T.; Jewell, K. S.; Dietrich, C.; Wick, A.; Ternes, T. A. Application of a non-
739 target workflow for the identification of specific contaminants using the example of the
740 Nidda river basin. *Water Research* **2020**, *178*, 115703.
- 741 (58) Kunzelmann, M.; Winter, M.; Åberg, M.; Hellenäs, K. E.; Rosén, J. Non-targeted anal-
742 ysis of unexpected food contaminants using LC-HRMS. *Anal. Bioanal. Chem.* **2018**,
743 *410*, 5593–5602.
- 744 (59) Reemtsma, T.; Berger, U.; Arp, H. P. H.; Gallard, H.; Knepper, T. P.; Neumann, M.;
745 Quintana, J. B.; Voogt, P. D. Mind the Gap: Persistent and Mobile Organic Compounds
746 - Water Contaminants That Slip Through. *Environ. Sci. Technol.* **2016**, *50*, 10308–
747 10315.
- 748 (60) Brüggem, S.; Schmitz, O. J. A New Concept for Regulatory Water Monitoring Via High-

- 749 Performance Liquid Chromatography Coupled to High-Resolution Mass Spectrometry.
750 *J Anal Test* **2018**, *2*, 342–351.
- 751 (61) Badea, S. L.; Geana, E. I.; Niculescu, V. C.; Ionete, R. E. Recent progresses in an-
752 alytical GC and LC mass spectrometric based-methods for the detection of emerging
753 chlorinated and brominated contaminants and their transformation products in aquatic
754 environment. *Sci. Total Environ.* **2020**, *722*, 137914.
- 755 (62) Greibrokk, T.; Andersen, T. High-temperature liquid chromatography. *J. Chromatogr.*
756 *A* **2003**, *1000*, 743–755.
- 757 (63) Snyder, L. R.; Kirkland, J. J.; Dolan, J. W. *Introduction to Modern Liquid Chromatog-*
758 *raphy*; John Wiley and Sons, 2010.
- 759 (64) Bade, R.; Bijlsma, L.; Sancho, J. V.; Hernández, F. Critical evaluation of a simple
760 retention time predictor based on LogKow as a complementary tool in the identification
761 of emerging contaminants in water. *Talanta* **2015**, *139*, 143–149.
- 762 (65) Kaliszan, R.; Haber, P.; Tomasz, B.; Siluk, D.; Valko, K. Lipophilicity and pKa esti-
763 mates from gradient high-performance liquid chromatography. *J. Chromatogr. A* **2002**,
764 *965*, 117–127.
- 765 (66) Phillips, A. L.; Peter, K. T.; Sobus, J. R.; Fisher, C. M.; Manzano, C. A.;
766 McEachran, A. D.; Williams, A. J.; Knolhoff, A. M.; Ulrich, E. M. Standardizing
767 non-targeted analysis reporting to advance exposure science and environmental epi-
768 demiology. *J. Expo. Sci. Environ. Epidemiol.* **2023**,
- 769 (67) Peter, K. T.; Phillips, A. L.; Knolhoff, A. M.; Gardinali, P. R.; Manzano, C. A.;
770 Miller, K. E.; Pristner, M.; Sabourin, L.; Sumarah, M. W.; Warth, B.; Sobus, J. R.
771 Nontargeted Analysis Study Reporting Tool: A Framework to Improve Research Trans-
772 parency and Reproducibility. *Anal. Chem.* **2021**, *93*, 13870–13879.

- 773 (68) Place, B. J. et al. An Introduction to the Benchmarking and Publications for Non-
774 Targeted Analysis Working Group. *Anal. Chem.* **2021**, *93*, 16289–16296.
- 775 (69) Malm, L.; Palm, E.; Souihi, A.; Plassmann, M.; Liigand, J.; Kruve, A. Guide to Semi-
776 Quantitative Non-Targeted Screening Using LC/ESI/HRMS. *Molecules* **2021**, *26*, 3524.
- 777 (70) Zubarev, R. A.; Makarov, A. Orbitrap mass spectrometry. *Anal. Chem.* **2013**, *85*,
778 5288–5296.
- 779 (71) Boesl, U. Time-of-flight mass spectrometry: Introduction to the basics. *Mass Spectrom.*
780 *Rev.* **2017**, *36*, 86–109.
- 781 (72) Gosetti, F.; Mazzucco, E.; Gennaro, M. C.; Marengo, E. Contaminants in water: non-
782 target UHPLC/MS analysis. *Environ. Chem. Lett.* **2016**, *14*, 51–65.
- 783 (73) Jeong, Y.; Da Silva, K. M.; Iturrospe, E.; Fujii, Y.; Boogaerts, T.; van Nuijs, A. L.;
784 Koelmel, J.; Covaci, A. Occurrence and contamination profile of legacy and emerging
785 per- and polyfluoroalkyl substances (PFAS) in Belgian wastewater using target, suspect
786 and non-target screening approaches. *J. Hazard. Mater.* **2022**, *437*, 129378.
- 787 (74) Xia, X.; Zheng, Y.; Tang, X.; Zhao, N.; Wang, B.; Lin, H.; Lin, Y. Nontarget Identifi-
788 cation of Novel Per- and Polyfluoroalkyl Substances in Cord Blood Samples. *Environ.*
789 *Sci. Technol.* **2022**, *56*, 17061–17069.
- 790 (75) Yu, N.; Wen, H.; Wang, X.; Yamazaki, E.; Taniyasu, S.; Yamashita, N.; Yu, H.; Wei, S.
791 Nontarget Discovery of Per- And Polyfluoroalkyl Substances in Atmospheric Particulate
792 Matter and Gaseous Phase Using Cryogenic Air Sampler. *Environ. Sci. Technol.* **2020**,
793 *54*, 3103–3113.
- 794 (76) Alygizakis, N. A.; Oswald, P.; Thomaidis, N. S.; Schymanski, E. L.; Aalizadeh, R.;
795 Schulze, T.; Oswaldova, M.; Slobodnik, J. NORMAN digital sample freezing platform:
796 A European virtual platform to exchange liquid chromatography high resolution-mass

- 797 spectrometry data and screen suspects in “digitally frozen” environmental samples.
798 *TrAC Trends in Analytical Chemistry* **2019**, *115*, 129–137.
- 799 (77) Rostkowski, P. et al. The strength in numbers: comprehensive characterization of house
800 dust using complementary mass spectrometric techniques. *Anal. Bioanal. Chem.* **2019**,
801 *411*, 1957–1977.
- 802 (78) Alygizakis, N. A.; Samanipour, S.; Hollender, J.; Ibáñez, M.; Kaserzon, S.; Kokkali, V.;
803 Van Leerdam, J. A.; Mueller, J. F.; Pijnappels, M.; Reid, M. J.; Schymanski, E. L.;
804 Slobodnik, J.; Thomaidis, N. S.; Thomas, K. V. Exploring the Potential of a Global
805 Emerging Contaminant Early Warning Network through the Use of Retrospective Sus-
806 pect Screening with High-Resolution Mass Spectrometry. *Environ. Sci. Technol.* **2018**,
807 *52*, 5135–5144.
- 808 (79) Podwojski, K.; Fritsch, A.; Chamrad, D. C.; Paul, W.; Sitek, B.; Stühler, K.; Mutzel, P.;
809 Stephan, C.; Meyer, H. E.; Urfer, W.; Ickstadt, K.; Rahnenföhrer, J. Retention time
810 alignment algorithms for LC/MS data must consider non-linear shifts. *Bioinformatics*
811 **2009**, *25*, 758–764.
- 812 (80) Clauser, K. R.; Baker, P.; Burlingame, A. L. Role of accurate mass measurement (± 10
813 ppm) in protein identification strategies employing MS or MS/MS and database search-
814 ing. *Anal. Chem.* **1999**, *71*, 2871–2882.
- 815 (81) Minkus, S.; Bieber, S.; Letzel, T. Spotlight on mass spectrometric non-target screening
816 analysis: Advanced data processing methods recently communicated for extracting,
817 prioritizing and quantifying features. *Anal. Sci. Adv.* **2022**, *3*, 103–112.
- 818 (82) Kleis, J. N.; Hess, C.; Germerott, T.; Roehrich, J. Sensitive Screening of New Psychoac-
819 tive Substances in Serum Using Liquid Chromatography-Quadrupole Time-of-Flight
820 Mass Spectrometry. *J. Anal. Toxicol.* **2022**, *46*, 592–599.

821 (83) Horai, H. et al. MassBank: a public repository for sharing mass spectral data for life
822 sciences. *J. Mass. Spectrom.* **2010**, *45*, 703–714.

823 **TOC Graphic**

824

