

Critical assessment of covered chemical space with 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), 20 Cornwall Street, Woolloongabba, QLD, 4102, Australia*

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

[§]*Contributed equally to this work*

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

Abstract

Non-targeted analysis (NTA) has emerged as a valuable approach for comprehensive monitoring of chemicals of emerging concern (CECs) in the exposome. The NTA approach, theoretically, is able to identify compounds with diverse physicochemical properties and sources. Non-targeted analysis methods, even though generic and wide scoping, have been shown to have limitations in terms of their coverage of the chemical space, as the number of the identified chemicals in each sample is very low (e.g. $\leq 5\%$). Investigating the chemical space covered by each NTA assay is crucial for understanding the limitations and challenges associated with the workflow from experimental methods to the data acquisition and data processing. In this review, we examined recent NTA studies published between 2017 and 2023 that employed liquid chromatography-high resolution mass spectrometry. The parameters used in each

13 study were documented and reported chemicals at the confidence level 1 and 2 were
14 retrieved. The chosen experimental setups and the quality of reporting were critically
15 evaluated and discussed. The findings revealed that only around 2% of the estimated
16 chemical space (i.e. Norman SusDat) was covered by the NTA studies investigated.
17 Little to no trend was found between the experimental setup and the observed cov-
18 erage, due to the generic and wide scope of NTA studies. The limited coverage of
19 chemical space by the NTA studies highlights the necessity for a more comprehensive
20 approach in experimental and data processing setups to enable the exploration of a
21 broader range of chemical space, with the ultimate goal of protecting human and envi-
22 ronmental health. Recommendations to further explore a wider range of the chemical
23 space were given.

24 **Synopsis**

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

27 **Introduction**

28 The chemical space of the human and environmental exposome is highly diverse and mostly
29 unknown^{1,2}. The chemical space generally refers to all possible organic structures present
30 in our surrounding environment³. Theoretical estimates of such structures have suggested
31 around 10^{60} unique structures with molecular weights less than 500 Da^{4,5}. This theoretical
32 chemical space incorporates both known and unknown unknowns^{3,6}. These chemicals can
33 cause adverse effects depending on their structures and the exposure levels. In fact, when
34 looking at the known unknowns (i.e. structures recorded in the chemical databases), several
35 of them have been shown to have adverse effects on the environmental and human health⁷⁻⁹.

36 Chemical prioritization has been one of the main means for dealing with the diversity

37 of chemical space in the human and environmental exposome^{1,10,11}. This consists of ex-
38 ploration of the literature for measured chemicals and their properties/toxicities as well
39 as national/international chemical registries¹². A combination of predicted properties and
40 toxicity is used to rank chemicals in the databases based on their potential impact on the
41 environment and human health¹³. Chemicals with a high potential of such impact are consid-
42 ered as chemicals of emerging concern (CECs)^{14,15}. To facilitate the chemical prioritization,
43 several databases consisting of chemical structures, the associated physicochemical proper-
44 ties (both measured and predicted), and their biological activities have been made publicly
45 available (e.g. PubChem, Norman Databases, and CompTox)^{12,16}. However, most of these
46 known unknowns remain unmeasured in environmental and biological matrices due to diffi-
47 culties associated with the inclusion of such large number of chemicals in routine monitoring
48 programs.^{9,11}

49 Non targeted analysis (NTA) combined with chromatography coupled with high resolu-
50 tion mass spectrometry (LC-HRMS) is considered as one of the most comprehensive methods
51 for the detection and identification of known and unknown unknowns in complex environ-
52 mental and biological samples^{17,18}. This approach utilizes a generic and wide scope strategy
53 for the sample preparation and analysis to maximize the coverage of the chemical space of
54 the sample^{3,11,17,19-27}. This typically results in very large and complex datasets (e.g. 5 GB
55 per sample) that must be pre-processed prior to the identification workflow^{27,28}. The NTA
56 data processing workflows include several steps from data conversion to library search and
57 the confidence assessment of the candidate spectra^{3,19,22-25}. Due to the complexity of such
58 datasets and sheer size of the chemical databases, the NTA workflows are not very sensitive
59 and do not result in a high percentage of identified chromatographic features^{29,30}. A more
60 sensitive but less comprehensive data processing alternative is suspect screening where the
61 chemicals of interest are known prior to the data processing workflow. This approach is more
62 sensitive in terms of limits of detection but is unable to detect unknown unknowns^{16,25,31}.
63 These two strategies are commonly employed together for the screening of complex environ-

64 mental and biological samples¹⁹.

65 The NTA strategy, even though powerful, has not been widely accepted within the reg-
66 ulatory framework due to reproducibility issues^{26,29,32}. Recent studies have indicated that
67 small changes in both experimental (e.g. data dependent vs data independent acquisition)
68 and data processing parameters may result in different outcomes and thus conclusions^{29,30}.
69 In fact, the aforementioned issues with NTA assays have sparked a debate in the scientific
70 community and have given start to a new wave of data processing tools development^{21,33}.
71 Additionally, several efforts have been put into better defining the much needed quality con-
72 trol and assurance for such experiments to be successful in detection and identification of
73 the known and unknown unknowns in complex environmental samples, thus better under-
74 standing the coverage of the analyzed chemical space¹⁹.

75 Several recently published reviews discuss in detail the impact of different steps on the
76 chemical space coverage through different experimental approaches^{3,19,22,23}. They cover both
77 data processing and experimental parameters including study scope, sampling and sample
78 treatment, instrumental conditions, data processing and treatment, and reporting. How-
79 ever, none of these reviews attempted to assess (i.e. quantify) the coverage of the chemical
80 space reached by the already conducted NTA environmental studies. Quantification of the
81 coverage of chemical space by an analytical method is not a trivial task. Theoretically, it
82 can be quantified as a number of identified compounds in the given sample divided by the
83 number of all compounds present in the chemical subspace of the sample. But practically,
84 this calculation is impossible, due to the complex chemical nature of samples and the number
85 of unknown constituents. Nevertheless, the investigation of experimentally explored chem-
86 ical space is highly relevant for the researchers to be aware of the limited coverage of the
87 associated chemical space.

88 In this review, we aim to quantify the coverage of the chemical space by recent environ-
89 mental studies and investigate the relationship between the selected experimental param-
90 eters and the explored chemical space. To quantify the covered chemical space via NTA,

91 we collected all recent studies that perform NTA (not suspect screening) and reported lev-
92 els 1 and 2³⁴, in terms of identification, structures. Additionally, we limited the scope of
93 this study to semi-polar and polar chemicals analyzable with liquid chromatography cou-
94 pled with high resolution mass spectrometry (LC-HRMS), resulting in a total of 57 pa-
95 pers. As an approximation of the chemical space the Norman SusDat database containing
96 around 60k unique chemicals with available PubChem CIDs (compound ID number) was
97 used (<https://www.norman-network.com/nds/susdat/susdatSearchShow.php>). We
98 collected a list of experimental and instrumental parameters, including sample preparation
99 (i.e. storage and extraction conditions), chromatographic separation (e.g. eluents, gradient
100 type, and injection volume), high resolution mass spectrometry settings (e.g. mass analyzer,
101 data acquisition mode, and polarity), and data processing workflows (e.g. mass and reten-
102 tion time tolerance, retention time domain alignment and databases used for the search).
103 We also noted any unreported parameters to identify the most commonly omitted settings.
104 Furthermore, we extracted information on the scope of the studies and samples analyzed.

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

111 **Methods**

112 **Selection of NTA studies**

113 This review is particularly focused on the development of the NTA approach in environmental
114 studies, specifically after the discussions regarding reproducibility were initiated³³. Thus,
115 we used the citation database Scopus to search for relevant studies published from 2017

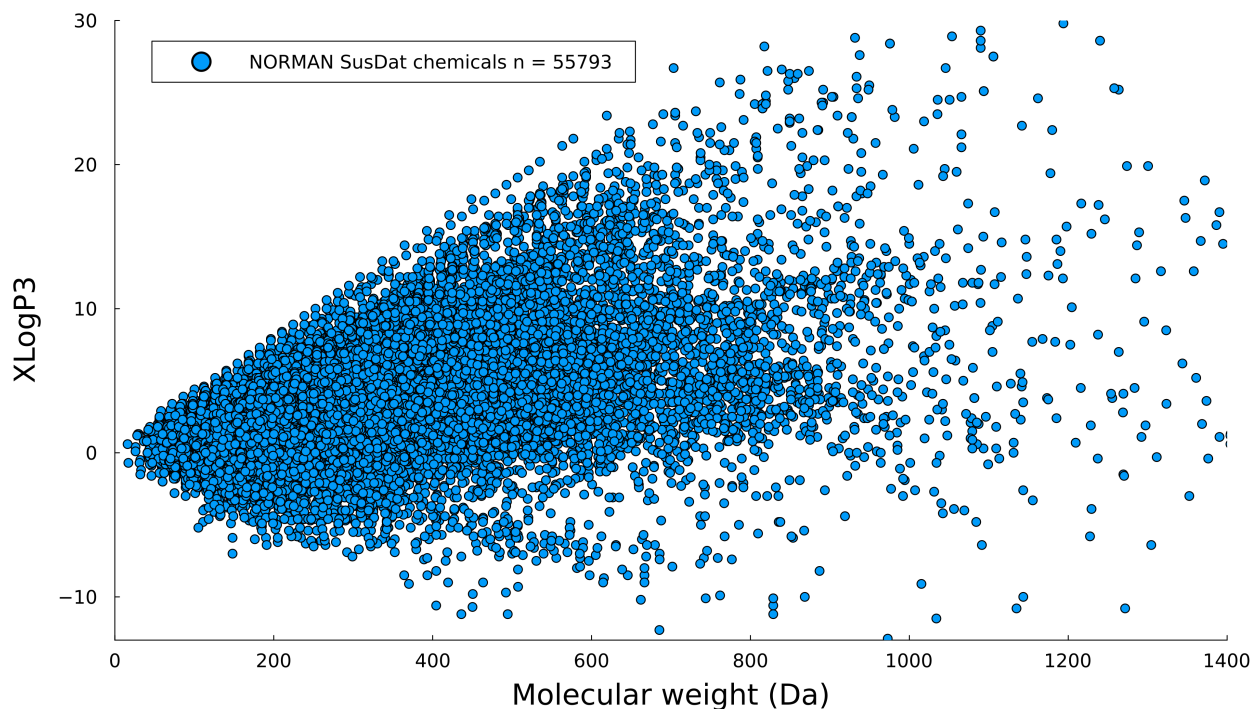


Figure 1: Distribution of all chemicals in the NORMAN SusDat database (n = 55793) based on their molecular weights (Da) and logP values.

116 to 2023 in the field of non-target analysis (NTA) with a focus on environmental science.
 117 The search was limited to articles that contained the keywords "non target analysis" or
 118 "untargeted analysis" or "untargeted screening" or "non-target screening" while excluding
 119 articles containing "metabolomics", "metabolic", or "gas chromatography". This initial
 120 search resulted in 377 publications adhering to the search parameters, which were then
 121 manually filtered to include only those that met a specific set of criteria.

122 The first criterion was that articles used non-target analysis to probe chemicals of emerg-
 123 ing concern, preferably in environmental matrices. Secondly, the publications had to use a
 124 non-target workflow. Some articles included the desired keywords in the title or abstract
 125 but were actually targeted studies with a very extensive list of target chemicals. The third
 126 criterion was that studies used LC-HRMS for sample analysis. Direct infusion studies, stud-
 127 ies that used rare setups, or heavily modified setups were excluded. Finally, review articles
 128 were excluded as they did not contribute any additional methods or identified compounds.
 129 The search for relevant studies meeting these criteria was completed on March 1st, 2023,

130 resulting in the inclusion of 57 studies in this review.

131 **Collection of instrumental parameters**

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

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

143 Where F is the flow rate (mL/min), $T\ run$ is the total run time of the method (min) -
144 excluding equilibration time- dc is the internal diameter of the column (cm), and L is the
145 length of the column in (cm). HRMS instrumental parameters included the mass analyzer,
146 sampling rate (in the case of Q-TOF), resolution (in the case of Orbitrap), data acquisition
147 mode, polarity, and mass range. Data processing parameters included mass tolerance, time
148 domain alignment, mass calibration, retention time tolerance, databases used, and total
149 database size (labeled as small if ≤ 1000 compounds or large if > 1000 compounds). A
150 summary of the collected parameters can be found in Figure 2. Furthermore, we made note of
151 parameters that were not reported in order to identify which settings were commonly omitted.

152 Lastly, we gathered information on the scope of the studies. The collected parameters along
153 with the list of the publications are publicly available through this link³⁵.

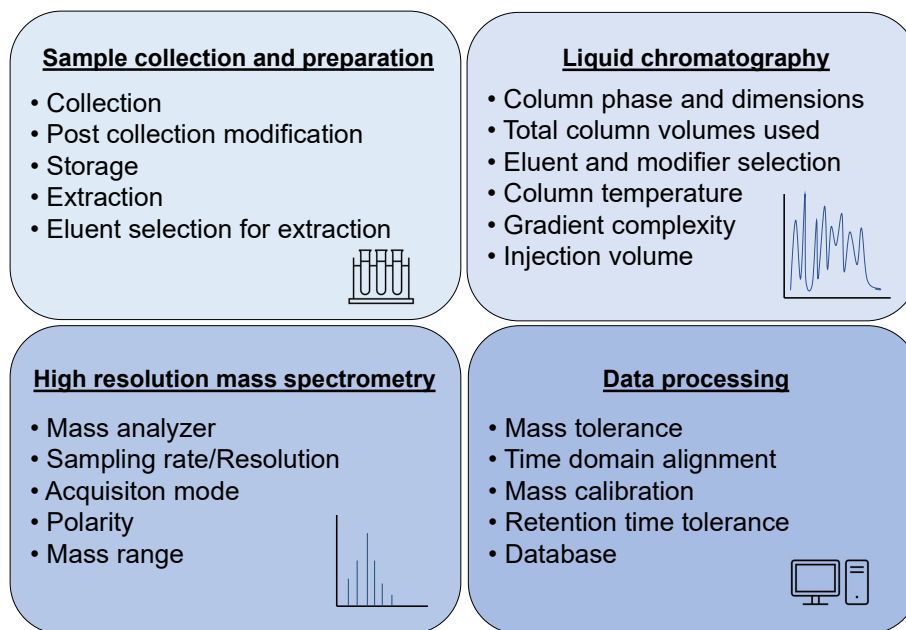


Figure 2: Summary of the main instrumental parameters collected from the reviewed NTA studies

154 Collection of reported structures

155 To assess the extent of chemical space coverage by recent NTA studies, we extracted the
156 reported structures. To ensure the reliability and accuracy of our analysis, we only included
157 structures identified with a high level of confidence (i.e levels one and two on the Schyman-
158 ski scale), which is less susceptible to false positive identifications³⁴. For each compound,
159 SMILES, IUPAC name, and the regular names provided by the authors were extracted. Fi-
160 nally, we excluded articles from our chemical space coverage assessment if the authors did
161 not specify the identification level, did not include the identified compounds in either the
162 article or supplementary materials, or only reported compounds within their target list.

163 Data processing

164 The list of the collected compounds was stored in CSV format, and Julia version 1.7 was
165 used to import and process the data. A modified version of the PubChemCrawler.jl package
166 was employed to retrieve chemical data such as XLogP3 and MW of the compounds from
167 the PubChem database by using their available identifiers (SMILES, IUPAC, InChIKey, or
168 regular name)³⁶. logP values extracted from PubChem are generated using XlogP3 with an
169 additive model starting from a reference compound³⁷. Retrieved data along with the col-
170 lected experimental parameters were combined into a dataset that included PubChem CIDs
171 corresponding to the compounds, their logP values, molecular weights, and experimental
172 parameters.

173 For the evaluation of the chemical space coverage, we additionally calculated elemental
174 mass defects (EMD) of six elemental ratios (CO, CCl, CN, CS, CF, and CH) for each
175 collected compound and the ones included in the NORMAN SusDat database³⁸. EMD
176 values were used to cluster structurally similar compounds together and separate others, as
177 they incorporate structural information and are used to compare compounds based on their
178 elemental composition³⁹. The combination of logP, MW, and EMDs was used for principal
179 component analysis (PCA), which is an unsupervised algorithm for dimensional reduction
180 combining variables into principal components⁴⁰. This approach is able to identify trends and
181 clusters in the data sets. Prior to the analysis the data was mean-centered and scaled to keep
182 the initial weight of all variables comparable. PCA was performed using the ScikitLearn.jl
183 julia package and in total three principal components were utilized.

184 The NORMAN SusDat database was used for the approximation of the chemical space
185 of environmental samples. While the chemical space comprises both known and unknown
186 compounds, it is practically impossible to include the latter in our approximations. The
187 Norman SusDat database includes CECs that have either been detected in various environ-
188 mental compartments or have been identified as potential CECs, providing a comprehensive
189 set of chemicals with a wide coverage of physical and chemical properties, and structures¹⁶.

190 Finally, the classes of the collected compounds were defined to illustrate the frequency
191 of identification of specific classes. To obtain the class of each CEC, the corresponding
192 InChIKey was used to generate information on superclasses, classes, and sub-classes of each
193 compound via ClassyFire. ClassyFire divides a given chemical compound into classes based
194 on its structural features (i.e. functional groups)⁴¹.

195 Discussion

196 In this review, we estimated the coverage of the chemical space of environmental samples by
197 investigating recent NTA studies. To evaluate the impact of selected workflow parameters
198 on the coverage of chemical space, we collected information on these parameters (e.g. mass
199 analyzer, data acquisition mode, ionization mode and size of the database used) from the
200 studies. The identified compounds were categorized into classes and their relative frequency
201 of occurrence was determined. XLogP3, MW, and EMDs were used to represent the vastness
202 of the chemical space, approximated with the NORMAN SusDat Database. PCA was em-
203 ployed to illustrate the coverage of the space of chemicals detected in recent environmental
204 studies.

205 Overview of the studies

206 In total, 57 studies were collected, with 54 of them published after 2019. Only studies using
207 NTA were included, while those using screening or targeted approaches but claiming to be
208 untargeted were excluded. Therefore, the significant increase in the number of such studies
209 in recent years reflects the successful development of NTA workflows. The scope of these
210 studies varies, with 27 studies focusing on a wide range of chemicals and another 20 studies
211 specifically targeting groups among which are per- and polyfluoroalkyl substances (PFAS),
212 pesticides, pharmaceuticals, and illicit drugs. Such prior prioritization influences the choice
213 of experimental setup. The remaining 10 studies focused on NTA workflow development,

214 indicating a growing interest and the need for further advancements in this field.

215 **Overview of selected parameters**

216 **Sample collection and preparation**

217 The collection and preparation of samples in the non-targeted analysis (NTA) workflow can
218 introduce potential sources of loss of chemical information. Issues such as ensuring sample
219 representativeness (e.g. selecting appropriate grab or passive sampling techniques), address-
220 ing potential sample contamination, accounting for matrix effects, optimizing extraction
221 methods for selectivity, and avoiding bias towards specific chemical groups are important
222 considerations in NTA^{3,19,22}. These challenges may impact the accuracy and reliability of
223 NTA results, potentially affecting the comprehensiveness and quality of the chemical infor-
224 mation obtained from the analysis. Therefore, careful attention to sample collection and
225 preparation steps are essential to minimize potential sources of bias and ensure robust and
226 reliable NTA outcomes.

227 The majority of the collected studies (57%) analyzed water samples (n = 38). Other
228 matrices that were investigated include biota (n = 5), dust (n = 3), urine (n = 3), atmospheric
229 particulate matter (n = 2), paper (n = 2), serum (n = 2), blood (n = 1), human hair (n =
230 1), ovarian follicular fluid (n = 1), sewage sludge (n = 1), snow (n = 1), and surface soil (n
231 = 1).

232 To prevent microbiological growth, the studies on water samples reported a conservation
233 step, which involved either adding an acid or storing the sample at a temperature of -20°C or
234 4°C. Out of the 38 water studies, 5 studies either did not include a step to stop microbiological
235 growth or did not report it. If this step was omitted, it could significantly alter the sample's
236 final composition when it is eventually analyzed in the laboratory^{42,43}.

237 Around 53% of publications analyzing water included a sample filtering step prior to
238 analysis. This step is a compromise to preserve the LC system and column but may lead
239 to the loss of the chemicals adsorbed to the particle's surface. Approximately 68% of stud-

240 ies included solid phase extraction (SPE) in their sample preparation, out of which 72%
241 used reversed-phase hydrophilic-lipophilic balance (HLB) SPE. However, only 33% of stud-
242 ies with SPE used acidic and/or basic modifiers in the extraction eluents. That implies that
243 most studies using only HLB SPE are potentially leaving ionizable compounds on the sor-
244 bent and may exclude them from the analysis. The remaining studies employed alternative
245 pre-treatment techniques among which vacuum-assisted evaporation, centrifugation, liquid-
246 liquid extraction (LLE), ultrasonic extraction as well as their combination. These choices are
247 mostly dictated by the sample nature/matrix. There were three studies that performed no
248 sample extraction and injected directly into the LC-MS with a higher injection volume⁴⁴⁻⁴⁶.
249 While this protocol minimizes sample adulteration and keeps the sampling of chemical space
250 more comprehensive, it can also pose a challenge to detection sensitivity due to the low
251 analyte concentration¹⁹.

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

257 **Liquid chromatography**

258 Chromatographic separation is employed to minimize sample complexity by spreading ana-
259 lytes across the time axis. This helps to reduce ion suppression (matrix effect) and provides
260 additional information (retention time) for the identification of the analytes. The chemistry
261 of the stationary phase along with the elution conditions affects the quality of separation
262 and the type of analytes being retained. Thus, the selection of chromatographic conditions
263 heavily influences the coverage of the chemical space of the sample¹⁹.

264 The majority of NTA studies use conventional reverse-phase separation with a generic
265 C18 column. Optimization of the separation includes proper selection of eluents and modi-

266 fiers, including suitable elution power and gradient setup, to avoid co-elution and excessive
267 or insufficient retention of chemicals⁴⁷. A simple linear gradient of an aqueous phase and
268 methanol or acetonitrile from low to high percentage is most widely accepted for the wide
269 scope screening. This method proved its reproducibility across different scopes of the stud-
270 ies²². However, this strategy focuses on polar to semipolar compounds, potentially excluding
271 very polar (i.e logP smaller than -2) and very hydrophobic substances (i.e. logP larger 6)
272 from the comprehensive investigation of the chemical composition of samples⁴⁸. To cover
273 the polar part of the chemical space, orthogonal methods such as hydrophilic interaction
274 chromatography (HILIC) become more popular while for hydrophobic volatile chemicals,
275 GC is a widely used technique^{44,49,50}. Finally, to ensure the reproducibility and reliability of
276 the studies parameters such as injection volume and column temperature should be properly
277 reported⁵¹.

278 More than 90% of the collected studies used a C18 column for the separation, among
279 which almost all were endcapped with a column length of 50mm (18%), 100mm (50%), or
280 150mm (31%). Column diameters were either 2.1mm (80% of the studies), 3mm (17%) with
281 the particle diameter under 3.5 μm . Additionally, two different studies reported 4.6mm and
282 0.05mm column diameters. Although applying a simple gradient ensures higher reproducibil-
283 ity of the method, only half of the studies (approximately 49%) used a linear gradient, while
284 around 32% used a semi-linear gradient and the remaining (18 %) used a more complex type
285 of gradient.

286 The median number of column volumes eluted in the studies is 16.2, with an interquartile
287 range of 15.6. The use of a sufficient number of column volumes should ensure the complete
288 elution of most hydrophobic compounds (high logP and MW) and the absence of carryover.
289 The optimal number depends on the stationary phase, eluent power, and analytes them-
290 selves⁵². Nevertheless, the widely accepted hypothesis is that there is a linear relationship
291 between logP and retention/number of column volumes used. The hypothesis is applied
292 for the reverse phase mode with comparable C18 selectivity, similar gradients, and eluent

293 composition^{53,54}. However, our results do not indicate the presence of a linear relationship
294 between the number of column volumes and logP of the chemicals, since no clear linear
295 pattern could be identified between these parameters (Figure S1).

296 In addition, the column temperatures used were all slightly above room temperature
297 which is favorable for repeatability and reproducibility⁵¹. 32% of publications used 40°C,
298 18% used 35°C, 13% used 30°C, two studies held the column at 25°C, one at 20°C, one at
299 45°C and one at 50°C. About 29% of papers did not report the column temperature, which
300 hinders the reproducibility of experiments.

301 Finally, 16% of the studies did not report the injection volume used. Injection volume
302 should not have a large effect on the final observed chemical space as they depend on the
303 extraction method and efficiencies. Nevertheless, the success of method's transfer depends
304 on it. The most studies used either 5 (n = 15) or 10 μ L (n = 12) injection volume, which is
305 adequate when using SPE extraction. The remaining were spread across 1, 3, 4, 7, 20, 100,
306 140, and 660 μ L.

307 To conclude, despite the rising discussion about reporting quality¹⁹, chromatographic
308 separation parameters in the collected studies were not always properly reported. Proper
309 harmonized reporting ensures successful method transfer, whereas inconsistent reporting
310 raises questions related to the reproducibility of the study, reliability of the results, and the
311 possibility of retrospective studies. While the majority of the studies seek to comprehensively
312 investigate the chemical composition of the samples, only approximately 10% employ an
313 alternative to the conventional approach to analyze the samples. Lastly, the hypothetical
314 linear trend between logP and retention was not confirmed, indicating the need for more
315 sophisticated strategies for the method development and optimization.

316 **High resolution mass spectrometry**

317 The Orbitrap and the quadrupole time of flight (QTOF) equipped with electrospray ioniza-
318 tion (ESI) are the two most commonly used HRMS instruments in liquid chromatography-

319 based (LC) NTA experiments. For complimentary analysis, it is recommended to perform
320 separate experiments in both positive and negative modes⁵⁵. The mass resolution of Orbitrap
321 mass analyzers is generally higher than that of QTOF, but both can provide high-resolution
322 mass spectra (*Resolution* $\geq 30,000$)⁵⁶.

323 In QTOF, resolution is determined by the architecture of the mass analyzer⁵⁷, while
324 for Orbitrap, the resolving power depends on a user specified resolution. In the case of
325 Orbitrap, the speed of scans is directly related to the spectral resolution. However, the
326 increase in mass resolution is limited by the time required for scanning operations. For
327 QTOF, a crucial parameter for data quality is the sampling speed, which is reported as
328 spectra per second in Hz. If the scan rate is too high, fewer ions are sampled, which can
329 lead to a sensitivity issue. Conversely, if the scan rate is too low, fewer data points on the
330 time axis are recorded, potentially causing missed detection of analytes eluting in a narrow
331 time range⁵⁸.

332 MS/MS spectra for structure elucidation are recorded using either data-dependent ac-
333 quisition (DDA) or data-independent acquisition (DIA). DDA mode records fragments of
334 pre-selected precursor ions, while DIA mode fragments all precursor ions within a certain
335 mass range. The latter is preferable for comprehensive investigations of complex samples.
336 However, DDA mode is currently the preferred choice in environmental studies, partly due
337 to the limited availability of processing tools for DIA files and also because the DIA experi-
338 mental setup is not commonly employed with Orbitrap mass analyzers²². QTOF analyzers
339 are more commonly used for DIA due to higher data acquisition rates.

340 Roughly, half of the collected studies ($n = 30$) utilized an Orbitrap mass analyzer, while
341 the other half ($n = 27$) employed a QTOF mass analyzer. However, a significant proportion
342 (approximately 74%) of the studies reported using DDA, which inherently limits their results
343 to predefined ions. The scan rate for QTOF analyzers was mostly set at 4 Hz, although
344 some studies operated at lower rates of 3, 2, or 1 Hz. Many studies using Orbitrap analyzers
345 operated at a resolution of 70,000, while some studies used lower resolutions with a minimum

346 of 35,000 and higher with a maximum of 240,000. Approximately 22% of the studies did not
347 report either resolution or scan rate.

348 Less than half of the studies (around 44%) conducted separate experiments in positive
349 and negative modes, utilizing multiple injections, different modifiers, and sometimes different
350 columns, which is considered a more suitable scenario for achieving comprehensive coverage of
351 chemical space. In approximately 30% of the studies, MS was operated only in positive mode.
352 There were ten publications where the analysis was reported in both modes, but the details
353 were insufficient to determine if the experiment was performed simultaneously or separately
354 in both modes. In three other studies, an exclusively negative mode was used to prioritize
355 a specific group of compounds of interest, such as PFAS⁵⁹⁻⁶¹, deliberately narrowing down
356 the investigated chemical space. Finally, two of the reviewed studies employed simultaneous
357 positive and negative ionization modes with formic acid as a modifier. This approach is
358 not preferable for NTA given that acidic additives are not always the optimal for a negative
359 ionization mode. Additionally, the acquired data becomes extremely complex and lacks
360 quality for reliable and robust processing.

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

365 To conclude, despite recent advancements in DIA technology, DDA remains the predom-
366 inant choice in the reviewed studies. However, the recommended approach for improved
367 reproducibility and reliability of NTA studies, and to enhance coverage of chemical space in
368 environmental and metabolomics research, is to acquire data in DIA mode for initial screen-
369 ing and then continue with DDA for individual feature identification. Finally, in terms of
370 reproducibility the lack of comprehensively reported information hinders method transfer
371 and therefore it warrants actions towards a harmonized reporting strategy.

372 Data processing

373 Data processing is considered a major bottleneck in NTA workflows. It refers to a series of
374 procedures that starts with the data conversion and ends with the feature identification¹⁹.
375 One of the steps for reliable processing is the mass calibration, either external or internal.
376 During this step the measured m/z values of known structures are compared against theoret-
377 ical m/z values. These shifts/correction factors are applied to all mass channels, depending
378 on the instrumental setup. This step ensures the quality of the spectra in terms of accu-
379 rate mass measurement⁶². An inadequate mass calibration may result in false positive and
380 negative detections during the identification⁶³.

381 One of the last steps of CEC identification is the use of a database to relate the MS out-
382 put to a known chemical structure. To proceed with the identification, experimental data
383 undergoes pre-processing steps: data compression, to remove noise and blank peaks, fea-
384 ture detection, to find features in 3-dimensional data, componentization, to group fragments
385 and isotopologues belonging to the same compound, and feature prioritization to reduce the
386 number of irrelevant features⁶⁴. Since most of the collected studies used vendor software for
387 the latter four steps, which makes it almost impossible to retrieve the information of algo-
388 rithms utilized, these parameters cannot be adequately discussed for their influence on the
389 coverage of chemical space. For the identification of known unknowns, pre-processed data is
390 compared with chemical databases and matched against references from available spectral
391 libraries, utilizing a combination of features, retention time, accurate mass, and fragmenta-
392 tion pattern³⁴. The mass and the retention tolerance are two initial parameters used for the
393 candidates' list compilation. These parameters, along with the database used, heavily affect
394 the results of the candidate search. The number of chemicals included in databases used in
395 the evaluated studies differs from a few hundred structures in in-house libraries⁶⁵ to tens and
396 hundreds of thousands in publicly available libraries²³ such as NORMAN,¹⁶ MassBank⁶⁶ or
397 PubChem⁶⁷. These search algorithms result in a set of candidate structures that ultimately
398 must be confirmed via either reference standard and/or an orthogonal method³⁴.

399 For the transparency and reproducibility of the method, proper reporting of applied se-
400 tups for each data processing step is essential. Nevertheless, a significant part of the studies
401 did not provide sufficient information to reproduce the results. Specifically, approximately
402 39% did not mention anything about mass calibration, while 26% reported that they per-
403 formed calibration but did not describe the procedure. Only about 35% included a report
404 on the mass calibration procedure. A large number of the papers (40%) also did not report
405 whether a retention alignment was performed. 35% did report the fact that a retention
406 alignment was done but did not specify the algorithm that was used or provided the details
407 on the parameters used. The remaining 25% of publications did report both the fact that
408 one was performed and which algorithm was used.

409 In contrast, mass tolerance applied for the search was reported in almost all studies,
410 around 95%. Among which around 86% used a mass tolerance for the database query of
411 5ppm, that is the unofficially accepted standard for the NTA database search. There were
412 also studies that used a relatively high mass tolerance of 17 ppm or 20ppm and some studies
413 that used mass tolerances lower than 5 at 3ppm, 2ppm, and even 1ppm. Generally, the
414 studies that were using the lower mass tolerances for the database search reported a higher
415 resolution of the mass analyzer. On the other hand, retention tolerances had much lower
416 reporting rates as 39% of the studies did not include this information. The remaining studies
417 used tolerances in a range between 0.1 min and 0.5 min. However, there are a few publications
418 that used a wider tolerance, up to 1.8 min, which may result in a high false positive rate.
419 Finally, approximately 12% did not report the databases used or referred to the software but
420 not the databases that the software was using. The majority, 81%, used a total database
421 size containing more than 5000 compounds, while only 5 studies used databases with less
422 than one thousand compounds.

423 The data processing step is one of the main bottlenecks for the NTA approach and thus
424 requires greater attention within the community. Nevertheless, the reporting quality needs
425 improvement. Furthermore, it was found that around 70% of the identified chemicals are

426 available in MassBank EU. That means that roughly 30% of the HRMS spectra acquired for
427 the identified compounds have not been deposited in public databases such as MassBank.
428 For NTA to reach its full potential, the expansion of publicly available spectral databases is
429 vital for the improvement of the coverage of chemical space at the identification step.

430 **Explored chemical space**

431 The studies yielded a total of 2277 compounds reported in the identification level 1 up to 2b,
432 of which 1416 are unique structures. However, in 7 studies, there was no report of either iden-
433 tification level, or any identifiers, which hinders the retrieval of the compounds from these
434 studies. The class of each collected CEC was obtained and displayed in Figure 3. The most
435 commonly found compounds were benzoids, followed by organocyclic compounds and then
436 organic acids and derivatives. The latter category, along with organohalogen compounds,
437 constitutes PFAS, which have been of particular interest in recent years. The median molec-
438 ular weights of compounds from SusDat were 239 Da and 257 Da for those collected from the
439 studies, with a median XLogP3 of 3.2 for SusDat and 2.8 for collected compounds. Based
440 on histograms in Figure 4, compounds with the most frequently occurring properties are
441 being identified in recent NTA studies, which can be partially explained by the generalized
442 experimental workflows with reverse phase C18 columns.

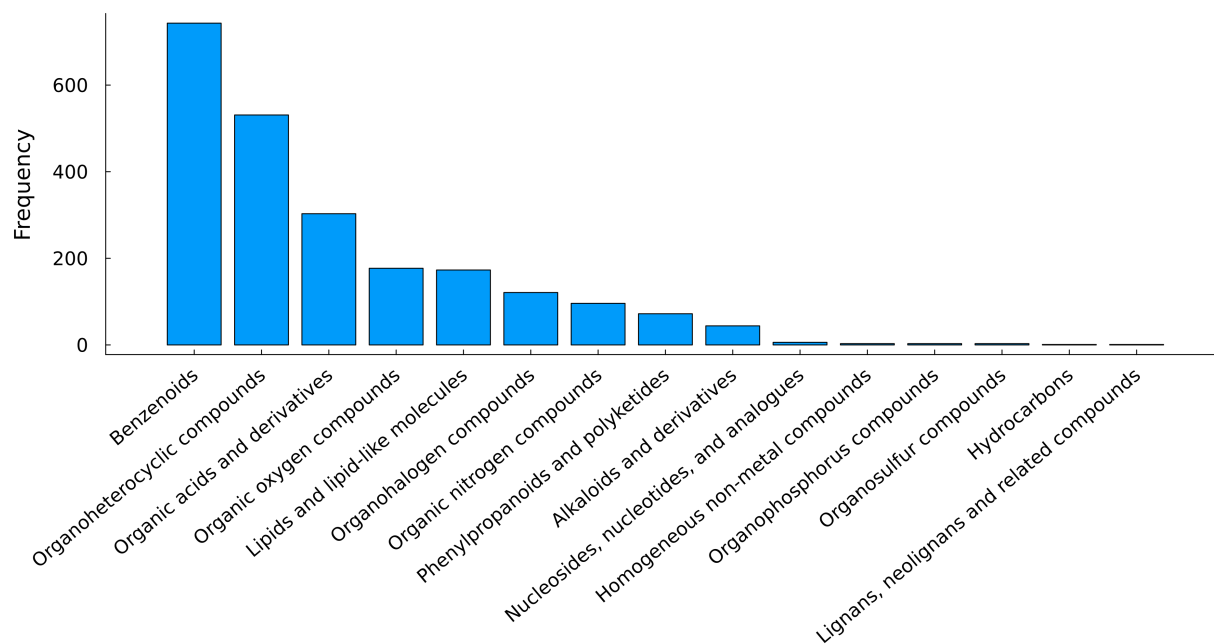


Figure 3: Histogram of all of the classes obtained from the Classyfire search for the detected CECs in reviewed studies

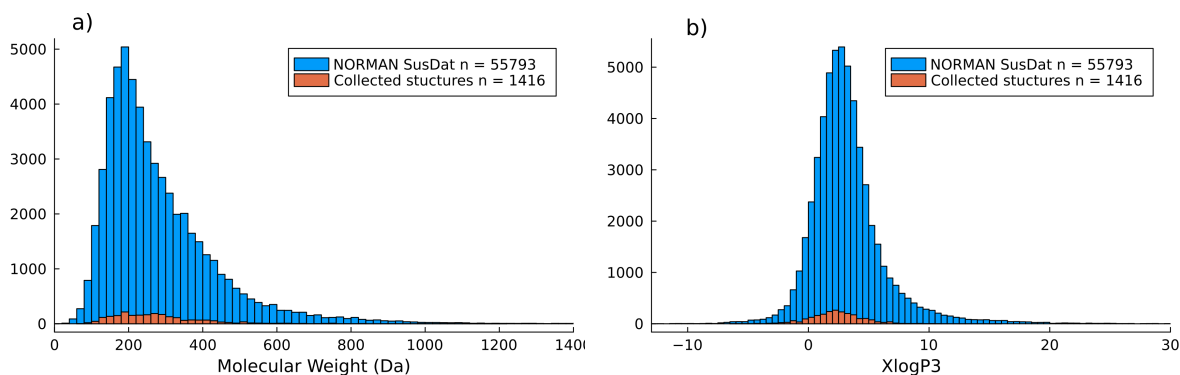


Figure 4: Molecular weights (a) and logP (b) distributions for the collected compounds (orange) and ones included in NORMAN susdat database (blue).

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

447 molecular weight (MW) and XLogP3, which emphasizes the limited space that is currently
448 explored using current non-target analysis (NTA) workflows. To examine the effect of some of
449 the mass spectrometry (MS) parameters used on the explored chemical space, all compounds
450 were plotted and clustered based on factors such as the mass analyzer used, acquisition mode,
451 ionization mode, and the total database size used (Figure S2-S5). However, neither of these
452 parameters showed an unambiguous influence on the coverage of the chemical space. It should
453 be noted that the representation in MW and logP dimensions does not provide information
454 about the elemental composition of compounds or their classes, which may result in an over-
455 representation of the covered chemical space. Therefore, it is important to consider other
456 parameters beyond MW and logP when evaluating the coverage of the chemical space by
457 the collected structures.

458 The PCA scores plot in Figure 6 reveals that many regions of the chemical space are
459 unexplored. The PCA was applied to the dataset combined the collected compounds with
460 the ones from the Norman SusDat, with MW, XLogP3, and the EMDs as input variables.
461 The first two principal components in the analysis were found to be primarily influenced by
462 the elemental mass differences (EMDs) associated with compounds containing chlorine (Cl),
463 fluorine (F), cyanide (CN), and sulfur (S). These EMDs represent the high variability in the
464 elemental composition of the compounds and were identified as the most important variables
465 in the PCA. This indicates that fewer compounds in the dataset contain halogens, nitrogen
466 (N), and sulfur, while hydrogen (H), which is present in every compound, does not contribute
467 significantly to the variability in the data. The third principal component is primarily
468 influenced by MW and XlogP3 (Figure S6). In total, the first three principal components
469 explain 74% of the variance (Figure S7). In Figure S8-S10, the coverage of chemical space
470 by different compound classes is displayed. Figure S10 specifically highlights the coverage by
471 organic acids and derivatives as well as organohalogen compounds. The majority of PFAS,
472 not exclusively, fall into these classes. The figures reveal that the distribution of compound
473 classes across the chemical space is not homogeneous, suggesting an over-representation of

474 certain classes of compounds. This observation can be attributed to the prior prioritization
475 of specific classes, which may bias the identification towards those classes of compounds.

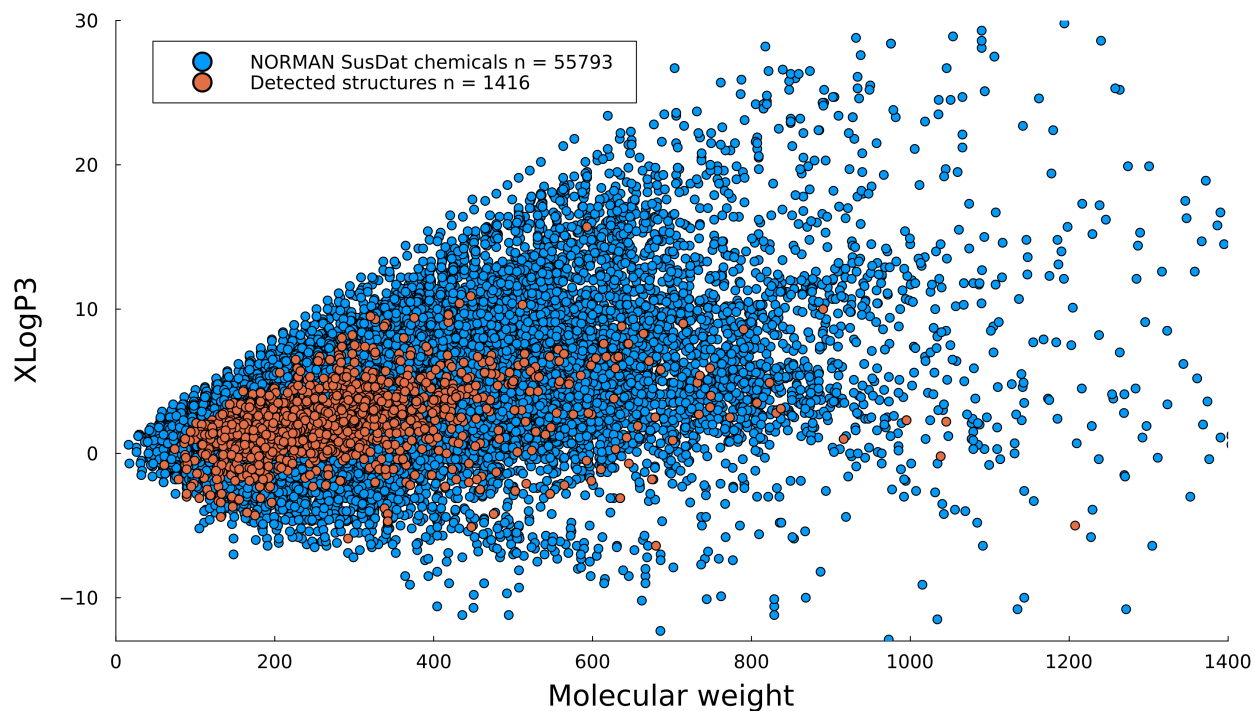


Figure 5: Distribution of all chemicals found in the reviewed articles at level 1 to 2b (orange) overlaid on NORMAN susdat database chemicals (blue) based on their molecular weights and XlogP3 value

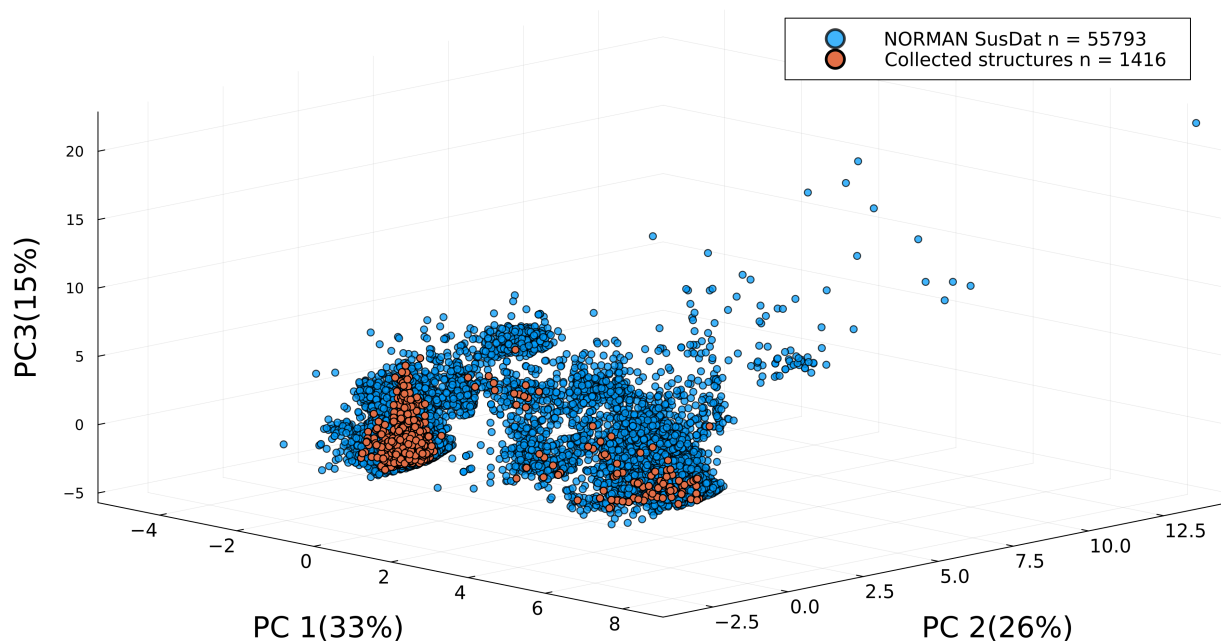


Figure 6: Scores plot of three principal components of the NORMAN susdat database (blue) and the collected structures (orange).

476 Overall only around 2% of the estimated chemical space was covered by NTA studies
 477 investigated in this review. This estimation is based on an approximation that by definition
 478 is far smaller than the true chemical space of the human and environmental exposome. No
 479 clear relationship between experimental conditions and coverage of the chemical space was
 480 discovered, which may indicate that the used experimental approaches are generic enough
 481 for the NTA assays. On the other hand, this may be caused by the lack of detailed and
 482 standardized reporting of the experimental conditions. Therefore, a more rigorous investi-
 483 gation of the parameters and standardization of reporting criteria has to be designed and
 484 performed. Although the most widely accepted properties of compounds such as logP and
 485 MW are widely used while discussing chemical space¹⁹, in this study we showed that they
 486 may not be the most relevant markers for assessing the coverage of chemical space. Finally,
 487 such a low coverage emphasizes the need for more comprehensive approaches to experimental
 488 and data processing workflows in order to explore a broader range of the chemical space and
 489 ultimately protect human and environmental health.

490 Recommendations and Outlook

491 Despite the ability of NTA to provide holistic information about the chemical composition of
492 the samples, their true coverage of the chemical space has not been investigated. Moreover,
493 the NTA studies have suffered from issues related to their reproducibility, due to the com-
494 plexity of both experimental and computational approaches employed in NTA assays. One
495 of the main bottlenecks for a more reproducible NTA assay is the lack of standardization of
496 the reporting criteria (including the experimental conditions). Our detailed investigation of
497 the previously published NTA studies further suggests the need for such criteria. Minimum
498 accepted experimental criteria and data processing parameters should be reported to ensure
499 the transparency and reliability of the results, which will potentially lead to the acceptance
500 of the NTA approach by the regulatory bodies.

501 The potential coverage of the chemical space should be assessed during the design of the
502 experimental setups. Most of the recent studies focused their experimental setups based
503 on the conventional workflow including HLB SPE for sample preparation, reverse phase
504 separation with C18 columns, and DDA acquisition mode, without considering alternative
505 approaches. The best practice would be an application of alternative extraction methods,
506 implementation of orthogonal techniques (e.g. RPLC and HILIC), DIA acquisition mode
507 as the first screening approach, and the application of reliable/robust data processing tools,
508 preferably open source/access. For the identification part of the workflow, the sharing of
509 experimental mass spectra of identified compounds is vital to the progress of the community.
510 Additionally, archiving the raw data in public repositories for both the retrospective analysis
511 as well as data processing tool development is highly essential.

512 To our knowledge, no other study has evaluated the coverage of the chemical space via
513 NTA studies in such detail. However, due to the lack of standardized reporting criteria,
514 the direct impact of different experimental choices on the covered chemical space could not
515 be established. Also, our study is limited to the works published after 2017 and we only
516 included studies with clear level 1 and 2 identification reporting. Moreover, we excluded the

517 suspect screening studies, which may result in an underestimation of the coverage of NTA
518 studies. However, our study, even though limited, clearly shows the shortcomings of the
519 current NTA practices and the need for further development in different areas - including
520 experimental setup.

521 **Acknowledgement**

522 The authors are thankful to the members of Environmental Modeling & Computational
523 Mass Spectrometry (www.emcms.info). S. Samanipour and V. Turkina are thankful to the
524 UvA Data Science Center and ChemistryNL TKI for the financial support (projects Edified
525 and SCOPE). The authors thank Denice van Herwerden for her help in setting up the
526 calculations of elemental mass defects. J. O'Brien is the recipient of a National Health and
527 Medical Research Council (NHMRC) Investigator Grant (2009209) funded by the Australian
528 Government. The Queensland Alliance for Environmental Health Sciences (QAEHS), The
529 University of Queensland (UQ), gratefully acknowledges the financial support of Queensland
530 Health.

531 **Notes**

532 Information retrieved in this study can be found at [https://doi.org/10.5281/zenodo](https://doi.org/10.5281/zenodo.7774345)
533 [.7774345](https://doi.org/10.5281/zenodo.7774345). References to the reviewed studies and collected experimental parameters are
534 at All experimental parameters.xlsx. The script to perform the calculations is available at
535 <https://github.com//tobihul//Code-for-Critical-assessment-of-covered-chemical-space-with-LC-HRMS-non-targeted-analysis>. PubChemCrawler package is
536 available at <https://github.com/JuliaHealth/PubChemCrawler.jl>.

Supporting Information Available

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

References

- (1) Wambaugh, J. F. et al. New approach methodologies for exposure science. *Curr. Opin. Toxicol.* **2019**, *15*, 76–92.
- (2) Escher, B. I.; Stapleton, H. M.; Schymanski, E. L. Tracking complex mixtures of chemicals in our changing environment. *Science* **2020**, *367*, 388–392.
- (3) Milman, B. L.; Zhurkovich, I. K. The chemical space for non-target analysis. *Trends Analyt Chem* **2017**, *97*, 179–187.
- (4) Van Deursen, R.; Reymond, J. L. Chemical space travel. *ChemMedChem* **2007**, *2*, 636–640.
- (5) Reymond, J. L. The Chemical Space Project. *Acc. Chem. Res.* **2015**, *48*, 722–730.
- (6) Reymond, J. L.; Awale, M. Exploring chemical space for drug discovery using the chemical universe database. *ACS Chem. Neurosci.* **2012**, *3*, 649–657.
- (7) Vermeulen, R.; Schymanski, E. L.; Barabási, A.-L.; Miller, G. W. The exposome and health: Where chemistry meets biology. *Science* **2020**, *367*, 392–396.
- (8) Landrigan, P. J. et al. The Lancet Commission on pollution and health. *Lancet* **2018**, *391*, 462–512.
- (9) Petrie, B.; Barden, R.; Kasprzyk-Hordern, B. A review on emerging contaminants in wastewaters and the environment: Current knowledge, understudied areas and recommendations for future monitoring. *Water Research* **2015**, *72*, 3–27.

- 560 (10) Samanipour, S.; O'Brien, J. W.; Reid, M. J.; Thomas, K. V.; Praetorius, A. From
561 Molecular Descriptors to Intrinsic Fish Toxicity of Chemicals: An Alternative Approach
562 to Chemical Prioritization. *Environ. Sci. Technol.* **2022**,
- 563 (11) Dulio, V.; van Bavel, B.; Brorström-Lundén, E.; Harmsen, J.; Hollender, J.;
564 Schlabach, M.; Slobodnik, J.; Thomas, K.; Koschorreck, J. Emerging pollutants in
565 the EU: 10 years of NORMAN in support of environmental policies and regulations.
566 *Environ Sci Eur* **2018**, *30*.
- 567 (12) Williams, A. J.; Grulke, C. M.; Edwards, J.; McEachran, A. D.; Mansouri, K.;
568 Baker, N. C.; Patlewicz, G.; Shah, I.; Wambaugh, J. F.; Judson, R. S.; Richard, A. M.
569 The CompTox Chemistry Dashboard: A community data resource for environmental
570 chemistry. *J Cheminform* **2017**, *9*.
- 571 (13) Tian, Z. et al. A ubiquitous tire rubber-derived chemical induces acute mortality in
572 coho salmon. *Science* **2021**, *371*, 185–189.
- 573 (14) Sauv e, S.; Desrosiers, M. A review of what is an emerging contaminant. *Chem. Cent.*
574 *J.* **2014**, *8*, 15.
- 575 (15) Maddela, N. R.; Ramakrishnan, B.; Kakarla, D.; Venkateswarlu, K.; Megharaj, M.
576 Major contaminants of emerging concern in soils: a perspective on potential health
577 risks. *RSC Advances* **2022**, *12*, 12396–12415.
- 578 (16) Mohammed Taha, H. et al. The NORMAN Suspect List Exchange (NORMAN-SLE):
579 facilitating European and worldwide collaboration on suspect screening in high resolu-
580 tion mass spectrometry. *Environ Sci Eur* **2022**, *34*.
- 581 (17) Aceña, J.; Stampachiachiere, S.; P erez, S.; Barcel o, D. Advances in liquid chromatogra-
582 phy - High-resolution mass spectrometry for quantitative and qualitative environmental
583 analysis. *Anal Bioanal Chem* **2015**, *407*, 6289–6299.

- 584 (18) Picó, Y.; Barceló, D. Transformation products of emerging contaminants in the envi-
585 ronment and high-resolution mass spectrometry: A new horizon. *Anal Bioanal Chem*
586 **2015**, *407*, 6257–6273.
- 587 (19) Schulze, B.; Jeon, Y.; Kaserzon, S.; Heffernan, A. L.; Dewapriya, P.; O'Brien, J.;
588 Gomez Ramos, M. J.; Ghorbani Gorji, S.; Mueller, J. F.; Thomas, K. V.; Samanipour, S.
589 An assessment of quality assurance/quality control efforts in high resolution mass spec-
590 trometry non-target workflows for analysis of environmental samples. *TrAC Trends in*
591 *Analytical Chemistry* **2020**, *133*, 116063.
- 592 (20) Muir, D. C.; Howard, P. H. Are there other persistent organic pollutants? A challenge
593 for environmental chemists. *Environ. Sci. Technol.* **2006**, *40*, 7157–7166.
- 594 (21) Samanipour, S.; Martin, J. W.; Lamoree, M. H.; Reid, M. J.; Thomas, K. V. Letter to
595 the Editor: Optimism for Nontarget Analysis in Environmental Chemistry. *Environ.*
596 *Sci. Technol.* **2019**, *53*, 5529–5530.
- 597 (22) Menger, F.; Gago-Ferrero, P.; Wiberg, K.; Ahrens, L. Wide-scope screening of polar
598 contaminants of concern in water: A critical review of liquid chromatography-high
599 resolution mass spectrometry-based strategies. *Trends Environ. Anal. Chem.* **2020**, *28*,
600 e00102.
- 601 (23) Hollender, J.; Schymanski, E. L.; Singer, H. P.; Ferguson, P. L. Nontarget Screening
602 with High Resolution Mass Spectrometry in the Environment: Ready to Go? *Environ.*
603 *Sci. Technol.* **2017**, *51*, 11505–11512.
- 604 (24) Zedda, M.; Zwiener, C. Is nontarget screening of emerging contaminants by LC-HRMS
605 successful? A plea for compound libraries and computer tools. *Anal. Bioanal. Chem.*
606 **2012**, *403*, 2493–2502.
- 607 (25) Krauss, M.; Singer, H.; Hollender, J. LC-high resolution MS in environmental analysis:

- 608 From target screening to the identification of unknowns. *Anal. Bioanal. Chem.* **2010**,
609 *397*, 943–951.
- 610 (26) McCord, J. P.; Groff, L. C.; Sobus, J. R. Quantitative non-targeted analysis: Bridging
611 the gap between contaminant discovery and risk characterization. *Environ. Int.* **2022**,
612 *158*, 107011.
- 613 (27) Schmidt, T. C. Recent trends in water analysis triggering future monitoring of organic
614 micropollutants. *Anal. Bioanal. Chem.* **2018**, *410*, 3933–3941.
- 615 (28) Samanipour, S.; Reid, M. J.; Thomas, K. V. Statistical Variable Selection: An Alterna-
616 tive Prioritization Strategy during the Nontarget Analysis of LC-HR-MS Data. *Anal.*
617 *Chem.* **2017**, *89*, 5585–5591.
- 618 (29) Hohrenk, L. L.; Itzel, F.; Baetz, N.; Tuerk, J.; Vosough, M.; Schmidt, T. C. Comparison
619 of Software Tools for Liquid Chromatography-High-Resolution Mass Spectrometry Data
620 Processing in Nontarget Screening of Environmental Samples. *Anal. Chem.* **2020**, *92*,
621 1898–1907.
- 622 (30) Samanipour, S.; Baz-Lomba, J. A.; Alygizakis, N. A.; Reid, M. J.; Thomaidis, N. S.;
623 Thomas, K. V. Two stage algorithm vs commonly used approaches for the suspect
624 screening of complex environmental samples analyzed via liquid chromatography high
625 resolution time of flight mass spectroscopy: A test study. *J. Chromatograph A* **2017**,
626 *1501*, 68–78.
- 627 (31) Hernández, F. et al. The role of analytical chemistry in exposure science: Focus on the
628 aquatic environment. *Chemosphere* **2019**, *222*, 564–583.
- 629 (32) Brack, W.; Hollender, J.; de Alda, M. L.; Müller, C.; Schulze, T.; Schymanski, E.;
630 Slobodnik, J.; Krauss, M. High-resolution mass spectrometry to complement monitor-
631 ing and track emerging chemicals and pollution trends in European water resources.
632 *Environ Sci Eur* **2019**, *31*, 62.

- 633 (33) Hites, R. A.; Jobst, K. J. Is Nontargeted Screening Reproducible? *Environ. Sci. Technol.* **2018**, *52*, 11975–11976.
634
- 635 (34) Schymanski, E. L.; Jeon, J.; Gulde, R.; Fenner, K.; Ruff, M.; Singer, H. P.; Hollender, J.
636 Identifying small molecules via high resolution mass spectrometry: Communicating
637 confidence. *Environ. Sci. Technol.* **2014**, *48*, 2097–2098.
- 638 (35) Hulleman, T.; Turkina, V.; O’Brien, J.; Chojnacka, A.; Thomas, K. V.; Saminopour, S.
639 Data files for: Critical assessment of covered chemical space with LC-HRMS non-
640 targeted analysis. **2022**,
- 641 (36) Hulleman, T. Code for: Critical assessment of covered chemical space with LC-HRMS
642 non-targeted analysis. 2023; [https://github.com/tobihul/Code-for-Critical-a](https://github.com/tobihul/Code-for-Critical-assessment-of-covered-chemical-space-with-LC-HRMS-non-targeted-analysis)
643 [ssessment-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).
- 644 (37) Cheng, T.; Zhao, Y.; Li, X.; Lin, F.; Xu, Y.; Zhang, X.; Li, Y.; Wang, R.; Lai, L.
645 Computation of octanol-water partition coefficients by guiding an additive model with
646 knowledge. *J. Chem. Inf. Model.* **2007**, *47*, 2140–2148.
- 647 (38) van Herwerden, D.; O’Brien, J. W.; Choi, P. M.; Thomas, K. V.; Schoenmakers, P. J.;
648 Samanipour, S. Naive Bayes classification model for isotopologue detection in LC-
649 HRMS data. *Chemom. Intell. Lab. Syst.* **2022**, *223*, 104515.
- 650 (39) Sleno, L. The use of mass defect in modern mass spectrometry. *J. Mass. Spectrom.*
651 **2012**, *47*, 226–236.
- 652 (40) Hastie, T.; Tibshirani, R.; Friedman, J. *The Elements of Statistical Learning Data*
653 *Mining, Inference, and Prediction*; Springer, 2009.
- 654 (41) Djoumbou Feunang, Y.; Eisner, R.; Knox, C.; Chepelev, L.; Hastings, J.; Owen, G.;
655 Fahy, E.; Steinbeck, C.; Subramanian, S.; Bolton, E.; Greiner, R.; Wishart, D. S. Classy-

- 656 Fire: automated chemical classification with a comprehensive, computable taxonomy.
657 *J Cheminform* **2016**, *8*, 61.
- 658 (42) Kruve, A. Semi-quantitative non-target analysis of water with liquid
659 chromatography/high-resolution mass spectrometry: How far are we? *Rapid*
660 *Commun. Mass Spectrom.* **2019**, *33*, 54–63.
- 661 (43) Lyytikäinen, M.; Kukkonen, J. V.; Lydy, M. J. Analysis of pesticides in water and
662 sediment under different storage conditions using gas chromatography. *Arch. Environ.*
663 *Contam. Toxicol.* **2003**, *44*, 437–444.
- 664 (44) Been, F.; Kruve, A.; Vughs, D.; Meekel, N.; Reus, A.; Zwartsen, A.; Wessel, A.; Fis-
665 cher, A.; ter Laak, T.; Brunner, A. M. Risk-based prioritization of suspects detected
666 in riverine water using complementary chromatographic techniques. *Water Res.* **2021**,
667 *204*, 117612.
- 668 (45) Hu, L. X.; Olaitan, O. J.; Li, Z.; Yang, Y. Y.; Chimezie, A.; Adepoju-Bello, A. A.;
669 Ying, G. G.; Chen, C. E. What is in Nigerian waters? Target and non-target screening
670 analysis for organic chemicals. *Chemosphere* **2021**, *284*, 131546.
- 671 (46) Köppe, T.; Jewell, K. S.; Dietrich, C.; Wick, A.; Ternes, T. A. Application of a non-
672 target workflow for the identification of specific contaminants using the example of the
673 Nidda river basin. *Water Research* **2020**, *178*, 115703.
- 674 (47) Kunzelmann, M.; Winter, M.; Åberg, M.; Hellenäs, K. E.; Rosén, J. Non-targeted anal-
675 ysis of unexpected food contaminants using LC-HRMS. *Anal. Bioanal. Chem.* **2018**,
676 *410*, 5593–5602.
- 677 (48) Reemtsma, T.; Berger, U.; Arp, H. P. H.; Gallard, H.; Knepper, T. P.; Neumann, M.;
678 Quintana, J. B.; Voogt, P. D. Mind the Gap: Persistent and Mobile Organic Compounds
679 - Water Contaminants That Slip Through. *Environ. Sci. Technol.* **2016**, *50*, 10308–
680 10315.

- 681 (49) Brügggen, S.; Schmitz, O. J. A New Concept for Regulatory Water Monitoring Via High-
682 Performance Liquid Chromatography Coupled to High-Resolution Mass Spectrometry.
683 *J Anal Test* **2018**, *2*, 342–351.
- 684 (50) Badea, S. L.; Geana, E. I.; Niculescu, V. C.; Ionete, R. E. Recent progresses in an-
685 alytical GC and LC mass spectrometric based-methods for the detection of emerging
686 chlorinated and brominated contaminants and their transformation products in aquatic
687 environment. *Sci. Total Environ.* **2020**, *722*, 137914.
- 688 (51) Greibrokk, T.; Andersen, T. High-temperature liquid chromatography. *J. Chromatogr.*
689 *A* **2003**, *1000*, 743–755.
- 690 (52) Snyder, L. R.; Kirkland, J. J.; Dolan, J. W. *Introduction to Modern Liquid Chromatog-*
691 *raphy*; John Wiley and Sons, 2010.
- 692 (53) Bade, R.; Bijlsma, L.; Sancho, J. V.; Hernández, F. Critical evaluation of a simple
693 retention time predictor based on LogKow as a complementary tool in the identification
694 of emerging contaminants in water. *Talanta* **2015**, *139*, 143–149.
- 695 (54) Kaliszan, R.; Haber, P.; Tomasz, B.; Siluk, D.; Valko, K. Lipophilicity and pKa esti-
696 mates from gradient high-performance liquid chromatography. *J. Chromatogr. A* **2002**,
697 *965*, 117–127.
- 698 (55) Malm, L.; Palm, E.; Souihi, A.; Plassmann, M.; Liigand, J.; Kruve, A. Guide to Semi-
699 Quantitative Non-Targeted Screening Using LC/ESI/HRMS. *Molecules* **2021**, *26*, 3524.
- 700 (56) Zubarev, R. A.; Makarov, A. Orbitrap mass spectrometry. *Anal. Chem.* **2013**, *85*,
701 5288–5296.
- 702 (57) Boesl, U. Time-of-flight mass spectrometry: Introduction to the basics. *Mass Spectrom.*
703 *Rev.* **2017**, *36*, 86–109.

- 704 (58) Gosetti, F.; Mazzucco, E.; Gennaro, M. C.; Marengo, E. Contaminants in water: non-
705 target UHPLC/MS analysis. *Environ. Chem. Lett.* **2016**, *14*, 51–65.
- 706 (59) Jeong, Y.; Da Silva, K. M.; Iturrospe, E.; Fujii, Y.; Boogaerts, T.; van Nuijs, A. L.;
707 Koelmel, J.; Covaci, A. Occurrence and contamination profile of legacy and emerging
708 per- and polyfluoroalkyl substances (PFAS) in Belgian wastewater using target, suspect
709 and non-target screening approaches. *J. Hazard. Mater.* **2022**, *437*, 129378.
- 710 (60) Xia, X.; Zheng, Y.; Tang, X.; Zhao, N.; Wang, B.; Lin, H.; Lin, Y. Nontarget Identifi-
711 cation of Novel Per- and Polyfluoroalkyl Substances in Cord Blood Samples. *Environ.*
712 *Sci. Technol.* **2022**, *56*, 17061–17069.
- 713 (61) Yu, N.; Wen, H.; Wang, X.; Yamazaki, E.; Taniyasu, S.; Yamashita, N.; Yu, H.; Wei, S.
714 Nontarget Discovery of Per- And Polyfluoroalkyl Substances in Atmospheric Particulate
715 Matter and Gaseous Phase Using Cryogenic Air Sampler. *Environ. Sci. Technol.* **2020**,
716 *54*, 3103–3113.
- 717 (62) Podwojski, K.; Fritsch, A.; Chamrad, D. C.; Paul, W.; Sitek, B.; Stühler, K.; Mutzel, P.;
718 Stephan, C.; Meyer, H. E.; Urfer, W.; Ickstadt, K.; Rahnenföhrer, J. Retention time
719 alignment algorithms for LC/MS data must consider non-linear shifts. *Bioinformatics*
720 **2009**, *25*, 758–764.
- 721 (63) Clauser, K. R.; Baker, P.; Burlingame, A. L. Role of accurate mass measurement (± 10
722 ppm) in protein identification strategies employing MS or MS/MS and database search-
723 ing. *Anal. Chem.* **1999**, *71*, 2871–2882.
- 724 (64) Minkus, S.; Bieber, S.; Letzel, T. Spotlight on mass spectrometric non-target screening
725 analysis: Advanced data processing methods recently communicated for extracting,
726 prioritizing and quantifying features. *Anal. Sci. Adv.* **2022**, *3*, 103–112.
- 727 (65) Kleis, J. N.; Hess, C.; Germerott, T.; Roehrich, J. Sensitive Screening of New Psychoac-

- 728 tive Substances in Serum Using Liquid Chromatography-Quadrupole Time-of-Flight
729 Mass Spectrometry. *J. Anal. Toxicol.* **2022**, *46*, 592–599.
- 730 (66) Horai, H. et al. MassBank: a public repository for sharing mass spectral data for life
731 sciences. *J. Mass. Spectrom.* **2010**, *45*, 703–714.
- 732 (67) Wang, Y.; Xiao, J.; Suzek, T. O.; Zhang, J.; Wang, J.; Bryant, S. H. PubChem: a
733 public information system for analyzing bioactivities of small molecules. *Nucleic Acids*
734 *Res.* **2009**, *37*, 623–633.

735 **TOC Graphic**

736

