

1 **patRoon: Open source software platform for environmental**
2 **mass spectrometry based non-target screening**

3 Rick Helmus^{a*}, Thomas L. ter Laak^{a,b}, Annemarie P. van Wezel^a, Pim de Voogt^a and Emma L.
4 Schymanski^c

5

6 ^a Institute for Biodiversity and Ecosystem Dynamics, University of Amsterdam, P.O. Box
7 94240, 1090 GE Amsterdam, The Netherlands

8 ^b KWR Water Research Institute, Chemical Water Quality and Health, P.O. Box 1072, 3430
9 BB Nieuwegein, The Netherlands

10 ^c Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, L-4367
11 Belvaux, Luxembourg.

12 * Corresponding author: r.helmus@uva.nl

13 **Abstract**

14 Mass spectrometry based non-target analysis is increasingly adopted in environmental
15 sciences to screen and identify numerous chemicals simultaneously in highly complex
16 samples. However, current data processing software either lack functionality for
17 environmental sciences, solve only part of the workflow, are not openly available and/or are
18 restricted in input data formats. In this paper we present *patRoon*, a new *R* based open-
19 source software platform, which provides comprehensive, fully tailored and straightforward
20 non-target analysis workflows. This platform makes the usage, evaluation and mixing of
21 well-tested algorithms seamless by harmonizing various commonly (primarily open)

22 software tools under a consistent interface. In addition, *patRoon* offers various functionality
23 and strategies to simplify and perform automated processing of complex (environmental)
24 data effectively. *patRoon* implements several effective optimization strategies to
25 significantly reduce computational times. The ability of *patRoon* to perform a
26 straightforward and effective non-target analysis was demonstrated with real-world
27 environmental samples, showing that *patRoon* makes comprehensive (environmental) non-
28 target analysis readily accessible to a wider community of researchers.

29 **Keywords**

30 High resolution mass spectrometry, compound identification, non-target analysis,
31 computational workflows

32 **Introduction**

33 Chemical analysis is widely applied in environmental sciences such as earth sciences,
34 biology, ecology and environmental chemistry, to study e.g. geomorphic processes,
35 (chemical) interaction between species or the occurrence, fate and effect of chemicals of
36 emerging concern in the environment. The environmental compartments investigated
37 include air, water, soil, sediment and biota, and exhibit a highly diverse chemical
38 composition and complexity. The number and quantities of chemicals encountered within
39 samples may span several orders of magnitude relative to each other. Therefore, chemical
40 analysis must discern compounds at ultra-trace levels, a requirement that can be largely met
41 with modern analytical instrumentation such as liquid or gas chromatography coupled with
42 mass spectrometry (LC-MS and GC-MS). The high sensitivity and selectivity of these

43 techniques enable accurate identification and quantification of chemicals in complex sample
44 materials.

45
46 Traditionally, a ‘target analysis’ approach is performed, where identification and
47 quantitation occur by comparing experimental data with reference standards. The need to
48 pre-select compounds of interest constrains the chemical scope of target analysis, and
49 hampers the analysis of chemicals with (partially) unknown identities such as transformation
50 products and contaminants of emerging concern (CEC). In addition, the need to acquire or
51 synthesize a large number of analytical standards may not be feasible for compounds with a
52 merely suspected presence. Recent technological advancements in chromatography and
53 high resolution MS (HRMS) allows detection and tentative identification of compounds
54 without the prior need of standards [1]. This ‘non-target’ analysis (NTA) approach is
55 increasingly adopted to perform simultaneous screening of up to thousands of chemicals in
56 the environment, such as finding new CEC [1–6], identifying chemical transformation
57 (by)products [7–12] and identification of toxicants in the environment [13–16].

58
59 Studies employing environmental NTA typically allow the detection of hundreds to
60 thousands of different chemicals [17, 18]. Effectively processing such data requires
61 workflows to automatically extract and prioritize NTA data, perform chemical identification
62 and assist in interpreting the complex resulting datasets. Currently available tools often
63 originate from other research domains such as life sciences and may lack functionality or
64 require extensive optimization before being suitable for environmental analysis. Examples
65 include handling chemicals with low sample-to-sample abundance, recognition of

66 halogenated compounds, usage of data sources with environmentally relevant substances,
67 or temporal and spatial trends. Furthermore, existing tools solve only part of the workflow,
68 generally use differing and incompatible data formats and employ different user interfaces.
69 Hence, the need to learn, combine, optimize and sometimes develop or adapt various
70 specialized software tools, and perform tedious transformation of datasets currently hinders
71 further adoption of NTA, especially in more routine settings lacking appropriate in-house
72 computational expertise.

73

74 An NTA workflow can be generalized as a four step process (Figure 1) [1]. Firstly, data from
75 LC or GC-HRMS is either acquired or retrieved retrospectively, and pre-treated for
76 subsequent analysis (Figure 1a). This pre-treatment may involve conversion to open data
77 formats (e.g. mzML [19] or mzXML [20]) to increase operability with open-source software,
78 re-calibration of mass spectra to improve accuracy and centroiding [21] or other raw data
79 reduction steps to conserve space such as trimming chromatographs or filtering mass scans
80 (e.g. with the functionality from the ProteoWizard suite [22]). Secondly (Figure 1b), features
81 with unique chromatographic and mass spectral properties (e.g. retention time, accurate
82 mass, signal intensity) are automatically extracted and features considered equivalent
83 across sample analyses are grouped to allow qualitative and (semi-) quantitative comparison
84 further down the workflow. Thirdly (Figure 1c), the feature dataset quality is refined, for
85 instance, via rule-based filters (e.g. minimum intensity and absence in sample blanks) and
86 grouping of features based on a defined relationship such as adducts or homologous series
87 (e.g. “componentization”). Further prioritization during this step of the workflow is often
88 required for efficient data analysis, for instance, based on chemical properties (e.g. mass

89 defect and isotopic pattern), suspected presence (i.e. “suspect screening”) or intensity
90 trends in time and/or space (e.g. reviewed in [1]). Finally (Figure 1d), prioritized features are
91 annotated, for instance by assigning chemical formulae or compounds from a chemical
92 database (e.g. *PubChem* [23] or *CompTox* [24]) based on the exact mass of the feature. The
93 resulting candidates are ranked by conformity with MS data, such as match with theoretical
94 isotopic pattern and *in silico* or library MS fragmentation spectra, and study-specific
95 metadata, such as number of scientific references and toxicity data [1, 25].

96

97 Various open and closed software tools are already available to implement (parts of) the
98 NTA workflow. Commercial software tools such as *MetaboScape* [26], *UNIFI* [27], *Compound*
99 *Discoverer* [28] and *ProGenesis QI* [29] provide a familiar and easy to use graphical user
100 interface, may contain instrument specific functionalities and optimizations and typically
101 come with support for their installation and usage. However, they are generally not open-
102 source or open-access and are often restricted to proprietary data formats. This leads to
103 difficulties in data sharing, as exact algorithm implementations and parameter choices are
104 hidden, while maintenance, auditing or code extension by other parties is often not
105 possible. Many open-source or open-access tools are available to process mass
106 spectrometry data (e.g. [30, 31] and summarized in Table 1). While many tools were
107 originally developed to process metabolomics and proteomics data, approaches such as
108 *XCMS* [32] and *MZmine* [33] have also been applied to environmental NTA studies [6, 34].
109 Many open tools are easily interfaced with the *R* statistical environment [35] (Table 1).
110 Leveraging this open scripting environment inherently allows defining highly flexible and
111 reproducible workflows and increases the accessibility of such workflows to a wider

112 audience as a result of the widespread usage of *R* in data sciences. Various open tools
113 overlap in functionality (Table 1), and are implemented with differing algorithms or
114 employing different data sources. As a consequence, tools may generate different results, as
115 has been shown when generating feature data [36–40] or performing structural annotations
116 [25, 41]. Thus, a flexible platform to combine and evaluate various algorithms that is
117 independent of closed MS vendor input data formats is desired in order to tailor an optimal
118 NTA workflow to the particular study types and methodological characteristics.

119

120 *Table 1. Overview of commonly used open-source or open-access software tools to implement NTA workflows.*

121 **<Table from end of this document should be placed here>**

122

123 Here, we present an *R* based open-source software platform called *patRoon* ('pattern' in
124 Dutch) providing comprehensive NTA data processing from HRMS data pre-treatment,
125 detection and grouping of features, through to molecular formula and compound
126 annotation. In *patRoon*, various (primarily open) tools commonly used for NTA data
127 processing are harmonized within a consistent and easy to use interface. In addition, new
128 functionality is implemented that simplify and improve NTA data processing, such as
129 automated chemical annotation, visualization and reporting of results, comparing and
130 combining results from different algorithms, and data reduction and prioritization
131 strategies. The architecture of *patRoon* is designed to be extendible in order to
132 accommodate for rapid developments in the NTA research field.

133 **Implementation**

134 The implementation section starts with an overview of the *patRoon* workflows. Subsequent
135 sections provide details on additional functionality implemented by *patRoon* which relate to
136 data processing, annotation, visualization and reporting. Finally, a detailed description is
137 given of the software architecture. *patRoon* is then demonstrated in the Results and
138 discussion section. The software tools and databases used for the implementation of
139 *patRoon* are summarized in Additional file 1.

140 ***Workflow in patRoon***

141 *patRoon* encompasses a comprehensive workflow for HRMS based NTA (Figure 2). All steps
142 within the workflow are optional and the order of execution is largely customizable. Some
143 steps depend on data from previous steps (blue arrows) or may alter or amend data from
144 each other (red arrows). The workflow commonly starts with pre-treatment (PT) of raw
145 HRMS data. Next, feature data is generated, which consists of finding features (FTS) in each
146 sample, an optional retention time alignment step, and then grouping into “feature groups”
147 (FG). FTS and FG may be preceded by automatic parameter optimization (PO), or followed
148 by suspect screening (SUS). The feature data may then finally be used for componentization
149 (CMT) and/or annotation steps, which involves generation of MS peak lists (MSPL), as well
150 as formula and compound annotations (FOR/COM). At any moment during the workflow,
151 the generated data may be inspected, visualized and treated by e.g. rule based filtering.
152 These operations are discussed in the next section.

153

154 Several commonly used open software tools, such as *OpenMS* [52], *XCMS* [32], *MetFrag* [48]
155 and *SIRIUS* [54–58], and closed software tools, such as *Bruker DataAnalysis* [61] (chosen due
156 to institutional needs), are interfaced to provide a choice between multiple algorithms for
157 each workflow step (Additional file 3: Table S1). Customization of the NTA workflow may be
158 achieved by freely selecting and mixing algorithms from different software tools. For
159 instance, a workflow that uses *XCMS* to group features allows that these features originate
160 from other algorithms than those supported by *XCMS* (e.g. those from *OpenMS*), a situation
161 that would require tedious data transformation when *XCMS* is used standalone.

162

163 To ease parameter selection over the various feature finding and grouping algorithms, an
164 automated feature optimization (FO) approach was adopted from the isotopologue
165 parameter optimization (IPO) *R* package [62], which employs design of experiments to
166 optimize LC-MS data processing parameters [63]. IPO was integrated in *patRoon*, and its
167 code base was extended to (a) apply to other feature finding and grouping algorithms
168 supported by *patRoon* (i.e. *XCMS*, *OpenMS* and *enviPick*), (b) support isotope detection with
169 *OpenMS*, (c) perform optimization of qualitative parameters and (d) provide a consistent
170 output format for easy inspection and visualization of optimization results.

171

172 In *patRoon*, componentization (CMT) refers to consolidating different (grouped) features
173 with a prescribed relationship, which is currently either based on (a) highly similar elution
174 profiles (i.e. retention time and peak shape), which are hypothesized to originate from the
175 same chemical compound (based on [53, 59]), (b) participation in the same homologous
176 series (based on [64]) or (c) the (normalized) intensity profiles across samples (based on [4,

177 5, 65]). Components obtained by approach (a) typically comprise adducts, isotopologues
178 and in-source fragments, and the supported algorithms in *patRoon* annotate these using
179 chemical rules. Approach (b) uses the *nontarget R* package [44] to calculate series from
180 aggregated feature data from replicates. The interpretation of homologous series between
181 replicates is assisted by merging series with overlapping features in cases where this will not
182 yield ambiguities to other series. If merging would cause ambiguities, instead links are
183 created that can then be explored interactively and visualized by a network graph generated
184 using the *igraph* [66] and *visNetwork* [67] *R* packages (see Additional file 2: Figure S1).

185
186 During the annotation step, molecular formulae and/or chemical compounds are
187 automatically assigned and ranked for all features or feature groups. The required MS peak
188 list (MSPL) input data are extracted from all MS analysis data files and subsequently pre-
189 processed, for instance, by averaging multiple spectra within the elution profile of the
190 feature and by removing mass peaks below user-defined thresholds. All compound
191 databases and ranking mechanisms supported by the underlying algorithms are supported
192 by *patRoon* and can be fully configured. Afterwards, formula and structural annotation data
193 may be combined to improve candidate ranking and manual interpretation of annotated
194 spectra. More details are outlined in the section “MS peak list retrieval, annotation and
195 candidate ranking”.

196 ***Data reduction, comparison and conversion***

197 Various rule-based filters are available for data-cleanup or study specific prioritization of all
198 data obtained through the workflow (see Table 2), and can be inverted to inspect the data

199 that would be removed (i.e. negation). To process feature data, multiple filters are often
200 applied, however, the order may influence the final result. For instance, when features were
201 first removed from blanks by an intensity filter, a subsequent blank filter will not properly
202 remove these features in actual samples. Similarly, a filter may need a re-run after another
203 to ensure complete data clean-up. To reduce the influence of order upon results, filters for
204 feature data are executed by default as follows:

- 205 1. an intensity pre-filter, to ensure good quality feature data for subsequent filters;
- 206 2. filters not affected by other filters, such as retention time and *m/z* range;
- 207 3. minimum replicate abundance, blank presence and 'regular' minimum intensity;
- 208 4. repetition of the replicate abundance filter (only if previous filters affected results);
- 209 5. other filters that are possibly influenced by prior steps, such as minimum abundance
210 in feature groups or sample analyses.

211 Note that the above scheme only applies to those filters requested by the user, and the user
212 can apply another order if desired.

213
214 Further data subsetting allows the user to freely select data of interest, for instance,
215 following a (statistical) prioritization approach performed by other tools. Similarly, features
216 that are unique or overlapping in different sample analyses may be isolated, which is a
217 straightforward but common prioritization technique for NTA studies that involve the
218 comparison of different types of samples.

219

220 *Table 2. Major rule-based filtering functionality implemented in patRoon.*

Filter functionality	Features		Annotation		Processing	
	FTS	FG	MSPL	FOR	COM	CMT
Intensity threshold	X	X	X			
Feature properties ¹	X	X				
Max intensity deviation across replicates		X				
Minimum intensity above blank		X				
Minimum size or abundance		X				X
Top most abundant/highest scoring			X	X	X	
Minimum scoring				X	X	
Annotation ²				X	X	X
Organic matter rules ³				X		

FTS: features; FG: feature groups; MSPL: MS peak lists; FOR: formulae; COM: compounds; CMT: components; (1) Retention time, chromatographic peak width, m/z and mass defect range; (2) e.g. adducts, isotopologues, formula composition, neutral loss; (3) expected formula composition based on [68–71].

221

222 Data from feature groups, components or annotations that are generated with different
 223 algorithms (or parameters thereof) can be compared to generate a consensus by only
 224 retaining data with (a) minimum overlap, (b) uniqueness or (c) by combining all results (only
 225 (c) is supported for data from components). Consensus data are useful to remove outliers,
 226 for inspection of algorithmic differences or for obtaining the maximum amount of data
 227 generated during the workflow. The consensus for formula and compound annotation data
 228 are generated by comparison of Hill-sorted formulae and the skeleton layer (first block) of
 229 the InChIKey chemical identifiers [72], respectively. For feature groups, where different
 230 algorithms may output deviating retention and/or mass properties, such a direct
 231 comparison is impossible. Instead, the dimensionality of feature groups is first reduced by
 232 averaging all feature data (i.e. retention times, m/z values and intensities) for each group.
 233 The collapsed groups have a similar data format as ‘regular’ features, where the compared
 234 objects represent the ‘sample analyses’. Subjection of this data to a feature grouping
 235 algorithm supported by *patRoon* (i.e. from *XCMS* or *OpenMS*) then allows straightforward

236 and reliable comparison of feature data from different algorithms, which is finally used to
237 generate the consensus.

238

239 Hierarchical clustering is utilized for componentization of features with similar intensity
240 profiles or to group chemically similar candidate structures of an annotated feature. The
241 latter “compound clustering” assists the interpretation of features with large numbers of
242 candidate structures (e.g. hundreds to thousands). This method utilizes chemical
243 fingerprinting and chemical similarity methods from the *rcdk* package [73] to cluster similar
244 structures, and subsequent visual inspection of the maximum common substructure then
245 allows assessment of common structural properties among candidates (methodology based
246 on [74]). Cluster assignment for both CMT and COM approaches is performed automatically
247 using the *dynamicTreeCut R* package [75]. However, clusters may be re-assigned manually
248 by the desired amount or tree height.

249

250 Several data conversion methods were implemented to allow interoperability with other
251 software tools. All workflow data types are easily converted to commonly used *R* data types
252 (e.g. *data.frame* or *list*), which allows further processing with other *R* packages.
253 Furthermore, feature data may be converted to and from native *XCMS* objects (i.e.
254 *xcmsSet* and *XCMSnExp*) or exported to comma-separated values (CSV) formats
255 compatible with *Bruker ProfileAnalysis* or *TASQ*, or *MZmine*.

256 ***MS peak list retrieval, annotation and candidate ranking***

257 Data for MS and MS/MS peak lists for a feature are collected from spectra recorded within
258 the chromatographic peak and averaged to improve mass accuracies and signal to noise
259 ratios. Next, peak lists for each feature group are assigned by averaging the mass and
260 intensity values from peak lists of the features in the group. Mass spectral averaging can be
261 customized via several data clean-up filters and a choice between different mass clustering
262 approaches, which allow a trade-off between computational speed and clustering accuracy.
263 By default, peak lists for MS/MS data are obtained from spectra that originate from
264 precursor masses within a certain tolerance of the feature mass. This tolerance in mass
265 search range is configurable to accommodate the precursor isolation window applied during
266 data acquisition. In addition, the precursor mass filter can be completely disabled to
267 accommodate data processing from data-independent MS/MS experiments, where all
268 precursor ions are fragmented simultaneously.

269

270 The formula annotation process is configurable to allow a tradeoff between accuracy and
271 calculation speeds. Candidates are assigned to each feature group, either directly by using
272 group averaged MS peak list data, or by a consensus from formula assignments to each
273 individual feature in the group. While the latter inherently consumes more time, it allows
274 removal of outlier candidates (e.g. false positives due to features with poor spectra).
275 Candidate ranking is improved by inclusion of MS/MS data in formula calculation (optional
276 for *GenForm* [47] and *DataAnalysis*).

277

278 Formula calculation with *GenForm* ranks formula candidates on isotopic match (amongst
279 others), where any other mass peaks will penalize scores. Since MS data of “real-world”
280 samples typically includes many other mass peaks (e.g. adducts, co-eluting features,
281 background ions), *patRoon* improves the scoring accuracy by automatic isolation of the
282 feature isotopic clusters prior to *GenForm* execution. A generic isolation algorithm was
283 developed, which makes no assumptions on elemental formula compositions and ion
284 charges, by applying various rules to isolate mass peaks that are likely part of the feature
285 isotopic cluster (see Additional file 2: Figure S2). These rules are configured to accommodate
286 various data and study types by default. Optimization is possible, for instance, to (a)
287 improve studies of natural or anthropogenic compounds by lowering or increasing mass
288 defect tolerances, respectively, (b) constrain cluster size and intensity ranges for low
289 molecular weight compounds or (c) adjust to expected instrumental performance such as
290 mass accuracy. Note that precursor isolation can be performed independently of formula
291 calculation, which may be useful for manual inspection of MS data.

292

293 Compound annotation is usually the most time and resource intensive process during the
294 non-target workflow. As such, instead of annotating individual features, compound
295 assignment occurs for the complete feature group. All compound databases supported by
296 the underlying algorithms, such as *PubChem* [23], *ChemSpider* [76] or *CompTox* [24] and
297 other local CSV files, as well as the scoring terms present in these databases, such as *in silico*
298 and spectral library MS/MS match, references in literature and presence in suspect lists, can
299 be utilized with *patRoon*. Default scorings supported by the selected algorithm/database or
300 sets thereof are easily selectable to simplify effective compound ranking. Furthermore,

301 formula annotation data may be incorporated in compound ranking, where a ‘formula
302 score’ is calculated for each candidate formula, which is proportional to its ranking in the
303 formula annotation data. Execution of unattended sessions is assisted by automatic restarts
304 after occurrence of timeouts or errors (e.g. due to network connectivity) and automatic
305 logging facilities.

306 ***Visualization, reporting and graphical interface***

307 In *patRoon*, visualization functionality is provided for feature and annotation data (e.g.
308 extracted ion chromatograms (EICs) and annotated spectra), to compare workflow data (i.e.
309 by means of Venn, chord and UpSet [77] diagrams, using the *VennDiagram* [78], *circlize* [79]
310 and *UpSetR* [80] *R* packages, respectively) and others such as plotting results from
311 automatic feature optimization experiments and hierarchical clustering data. Reports can be
312 generated in a common CSV text format or in a graphical format via export to a portable
313 document file (PDF) or hypertext markup language (HTML) format. The latter are generated
314 with the *R Markdown* [81, 82] and *flexdashboard* [83] *R* packages, and provide an easy to
315 use interface for interactive sorting, searching and browsing reported data. As plotting and
316 reporting functionalities can be performed at any stage during the workflow, the data that is
317 included in the reports is fully configurable.

318

319 While *patRoon* is primarily interfaced through *R*, several graphical user interface tools are
320 provided to assist the (novice) user. Most importantly, *patRoon* provides a *Shiny* [84] based
321 tool that automatically generates a commented template *R* script from user input, such as
322 selection of MS data file input, workflow algorithms and other common workflow

323 parameters (Figure 3a). Secondly, chromatographic data of features may be inspected either
324 by automatic addition of EICs in a *Bruker DataAnalysis* session or with a *Shiny* based
325 interface (Figure 3b).

326 ***Software architecture***

327 *patRoon* is distributed as an *R* package. Its source code is primarily written in the *R*
328 language, with some support code written in C++ and JavaScript. Both *Microsoft Windows*
329 (hereafter referred to as *Windows*) and *Linux* platforms are supported (support for macOS is
330 envisaged in the future). Several external dependencies are required; notable examples are
331 in Additional file 3: Table S1. *GenForm* is automatically compiled during package installation.
332 For *Windows* platforms, an installation script is provided to install and configure *patRoon*
333 and all of its dependencies automatically. Documentation includes a handbook, tutorial and
334 full reference manual [85–88], which are produced with the *bookdown* [89, 90], *R*
335 *Markdown* and *roxygen2* [91] *R* packages, respectively. Example data is contained in the
336 *patRoonData R* package [92, 93].

337
338 An important design goal was to provide a consistent, generic and easy to use interface that
339 does not require the user to know the implementation and interfacing details of the
340 supported algorithms. Each workflow step is executed by a generator function that takes
341 the desired algorithm and its parameters as input and returns objects from a common set of
342 data formats (see Figure 4). Names for commonly used parameters supported by multiple
343 algorithms are standardized for consistency and defaults are set where reasonable.
344 Furthermore, the format of input data such as retention time units as well as formula and

345 adduct specifications are harmonized and automatically converted to the format expected
346 by the algorithm. Nearly all parameters from the underlying algorithm can be set by the
347 user, hence, full configurability of the workflow is retained wherever possible. Generic
348 naming schemes are applied to output data, which assist the user in comparing results
349 originating from different algorithms. All exported functions from *patRoon* verify user input
350 with the *checkmate* [94] package, which efficiently performs tests such as correctness of
351 value range and type, and prints descriptive messages if input is incorrect.

352

353 A set of generic methods are defined for workflow classes that perform general data
354 inspection, selection, conversion and visualization, irrespective of the algorithm that was
355 used to generate the object (see Table 3). Consequently, the implementation of common
356 function names for multiple output classes allows a predictable and consistent user
357 interface.

358

359 *Table 3. Common generic methods defined in patRoon to process workflow data.*

Generic	Purpose
<code>length()</code> , <code>show()</code> , <code>algorithm()</code> , <code>names()</code> , <code>groupNames()</code>	obtain general object information such as object length and unique identifiers for contained results
<code>filter()</code>	rule-based filtering operations
<code>[, [[, \$ operators</code>	subsetting or extracting data
<code>as.data.table()</code> , <code>as.data.frame()</code>	conversion to <code>data.table</code> or <code>data.frame</code> object
<code>unique()</code> , <code>overlap()</code>	extract unique or overlapping features across replicates
<code>consensus()</code>	generates a consensus between different objects of the same class
<code>plot()</code> , <code>plotEIC()</code> , <code>plotSpec()</code>	plot general, chromatographic and annotation data
<code>plotChord()</code> , <code>plotUpSet()</code> ,	comparison of feature data or workflow objects from different algorithms by chord, UpSet and Venn diagrams

plotVenn()

360
361 Several optimization strategies are employed in *patRoon* to reduce computational
362 requirements and times. Firstly, external command line (CLI) tools are executed in parallel to
363 reduce overall execution times for repetitive (e.g. per sample analysis or per feature)
364 calculations. Commands are queued (first in, first out) and their execution is handled with
365 the *processx* package [95]. Secondly, functions employing time intensive algorithms
366 automatically cache their (partial) results in a local *SQLite* database file, which is accessed
367 via the *DBI* [96] and *RSQLite* [97] *R* packages. Thirdly, performance critical code dealing with
368 *OpenMS* data files and loading chromatographic data was written in C++ (interfaced with
369 *Rcpp* [98–100]) to significantly reduce times needed to read or write data. Fourthly, the
370 output files from *OpenMS* tools are loaded in chunks using the *pugixml* software library
371 [101] to ensure a low memory footprint. Finally, reading, writing and processing (large)
372 internal tabular data is performed with the *data.table* *R* package, which is a generally faster
373 and more memory efficient drop-in replacement to the native tabular data format of *R*
374 (*data.frame*), especially for large datasets [102].

375
376 Interfacing with *ProteoWizard* [22], *OpenMS*, *GenForm*, *SIRIUS* and *MetFrag* occurs by
377 wrapper code that automatically executes the CLI tools and perform the data conversions
378 necessary for input and output files. An alternative interface to *MetFrag* is also provided by
379 employing the *metFrag* *R* package [103], however, in our experience this option is currently
380 significantly slower than the CLI and therefore not used by default. For tools that are not
381 readily controllable from *R* (i.e. *ProfileAnalysis*, *TASQ* and *MZmine*), interfacing occurs via
382 importing or exporting CSV files (only export is supported for *MZmine*). Finally, the

383 *RDCOMClient R* package [104] is used to interface with *Bruker DataAnalysis* via the
384 distributed component object model, which allows automation of *DataAnalysis* functionality
385 from *R* that otherwise would only be available via its integrated visual basic scripting
386 environment.

387
388 A continuous integration pipeline performs automated tests during development and
389 delivers files to simplify installation of *patRoon* and all its dependencies (Additional file 2:
390 Figure S3). More than 900 unit tests are performed (>80% code coverage) with the *testthat*
391 and *vdiffr R* packages [105, 106]. After successful test completion, binary *R* packages
392 (*Windows*) and *Docker* images (*Linux*) are generated to facilitate installation of *patRoon* with
393 tested and compatible dependencies.

394 **Results and discussion**

395 ***Benchmark and demonstration data***

396 The data used to benchmark and demonstrate *patRoon* were obtained with an LC-HRMS
397 analysis of two different influent and effluent samples from a drinking water treatment pilot
398 installation and a procedural blank (all in triplicate). The samples originate from an
399 experiment where a set of 18 common environmental contaminants (yielding 20 individual
400 chromatographic peaks, see Additional file 3: Table S2) were spiked prior to drinking water
401 treatment. The analyses were performed using an LC-HRMS Orbitrap Fusion system
402 (ThermoFisher Scientific, Bremen, Germany) operating with positive electrospray ionization.
403 Further analytical conditions are as described in [11]. The raw data files can be obtained
404 from [107].

405 ***Parallelization benchmarks***

406 Several benchmarks were performed to test the multiprocessing functionality of *patRoon*.
407 Tests were performed on a personal computer equipped with an Intel® Core™ i7-8700K CPU
408 (6 cores, 12 threads), 32 gigabyte RAM, SATA SSD storage and the *Windows 10* Enterprise
409 operating system. Benchmarks were performed in triplicate using the *microbenchmark R*
410 package [108]. Standard deviations were below ten percent (see Figure 5a). Benchmarking
411 was performed on *msConvert* (MC), *FeatureFinderMetabo* (FFM), *GenForm* (GF), *SIRIUS* (SIR)
412 and *MetFrag* (MF). The multiprocessing functionality was compared to native
413 multithreading for the tools that supported this (FFM, SIR and MF). In addition, the
414 performance of batch calculations with multiprocessing was compared with native batch
415 calculation modes of tools where possible (MC and SIR). Parallelization methods were tested
416 with 1-12 parallel processes or threads (i.e. up to full utilization of both CPU threads of each
417 core). Input conditions were chosen to simulate “simple” and “complex” workflows, where
418 the latter resulted in more demanding calculations with ~2-10x longer mean execution
419 times (Table 4). The caching functionality of *patRoon* was disabled, where appropriate, to
420 obtain representative and reproducible test results. Prior to benchmarking, candidate
421 chemical compounds from PubChem for MF tests were cached in a local database to
422 exclude influences from network connectivity. Similarly, general spectral data required to
423 post-process FFM results were cached, as this is usually loaded once during a workflow,
424 even with varying input parameters. The input features for GF tests that resulted in very
425 long individual run times (i.e. >30 seconds) were removed to avoid excessive benchmark
426 runtimes. Generating feature and MS peaklist input data for annotation related tests was
427 performed with *patRoon* using algorithms from *OpenMS* and *mzR* [109], respectively. Pre-

428 treatment of feature data consisted of removal of features with low intensity and lacking
 429 MS/MS data. The number of features for SIR (except tests with native batch mode) and MF
 430 benchmarks were further reduced by application of blank, replicate and intensity filters to
 431 avoid long total runtimes due to their relatively high individual run times. Finally, the feature
 432 dataset was split in low (0-500) and high (500-1000) *m/z* portions, which were purposed for
 433 execution of “simple” and “complex” experiments, respectively. For more details of the
 434 workflow and input parameters see the *R* script code in Additional file 4. The software tools
 435 used for benchmarking are summarized in Additional file 1.

436

437 *Table 4. Utilized conditions for “simple” (S) and “complex” (C) tests.*

	Test	Input conditions ¹	Executions	Mean individual run time ² (s)
<i>msConvert (MC)</i>	MC-S	Conversion centroided input	15	4.8
	MC-C	Centroiding and conversion non-centroided input	15	8.5
<i>FeatureFinderMetabo (FFM)³</i>	FFM-S	High intensity threshold	15	4.1
	FFM-C	Low intensity threshold	15	38
<i>GenForm (GF)</i>	GF-S	CHNO elements, low <i>m/z</i>	512	0.2
	GF-C	CHNOPS elements, high <i>m/z</i>	128	1.7
<i>SIRIUS (SIR)³</i>	SIR-S	CHNO elements, low <i>m/z</i>	152 (512 ⁴)	2.3
	SIR-C	CHNOPS elements, high <i>m/z</i>	44 (128 ⁴)	7.7
<i>MetFrag (MF)³</i>	MF-S	Limited scoring, narrow mass search (5 ppm), low <i>m/z</i> .	152	3.0
	MF-C	Thorough scoring, wide mass search (20 ppm), high <i>m/z</i> .	44	8.6

(1): Features with *m/z* 0 – 500 (low) and *m/z* 500 – 1000 (high); (2): based on a test run without parallelization (n=3); (3) supports (configurable) native multithreading; (4) number of executions for native batch mode benchmarks.

438 When multiprocessing was used all tests (except GF-S, discussed below) showed a clear
 439 downward trend in execution times (down to ~200%-500%), and optimum conditions were

440 generally reached when the number of parallel processes equaled the number of physical
441 cores (six, see Figure 5a). When algorithms are fully parallelized, execution times are
442 expected to follow an inverse relationship with the number of parallel process (i.e. $1/n$) and
443 this was observed most closely with MC, whereas execution times for other tools show a
444 less steep reduction. Furthermore, utilizing multiple threads per core (i.e. hyperthreading)
445 did not reduce execution times further and even slowed down in some cases (e.g. MF-C).
446 These deviations in scalability were not investigated in detail. Since they were more
447 noticeable under complex conditions, it is expected that this may be caused by (a) more
448 involved post-processing results after each execution, which is currently not parallelized,
449 and (b) increased memory usage, which may raise the overhead of context switches
450 performed by the operating system. Nevertheless, the experiments performed here clearly
451 show that the multiprocessing functionality of *patRoan* can significantly reduce execution
452 times of various steps in an NTA workflow.

453

454 An exception, however, was the test performed with *GenForm* with simple conditions (GF-
455 S), which exhibited no significant change in execution times with multiprocessing (Figure
456 5a). Due to the particularly small mean run times (0.2 seconds) of this test, it was
457 hypothesized that the overhead of instantiating a new process from *R* (inherently not
458 parallelized) dominated the overall run times. To mitigate this, a 'batch mode' was
459 implemented, where such process initiation occurs from a command shell sub-process
460 instead. Here, multiple commands are executed by the sub-process in series, and the
461 desired degree of parallelization is then achieved by launching several of these sub-
462 processes and evenly dividing commands amongst them. The maximum size of each series

463 (or “batch size”) is configurable, and represents a balance between reduction of process
464 initiation overhead and potential loss of effectively load balancing of, for instance,
465 commands with highly deviating execution times. Next, various batch sizes were tested for
466 GF, both with and without multiprocessing parallelization (Additional file 2: Figure S4). For
467 GF-S, execution times clearly decreased with increasing batch sizes, however, no further
468 reduction was observed with parallelism. In contrast, serial execution of GF-C was not
469 affected by varying batch size, whereas added parallelism reduced execution times for small
470 batch sizes (≤ 8), but significantly increased such times for larger sizes. The results
471 demonstrate that the typical short lived GF executions clearly benefit from batch mode. In
472 addition, it is expected that by further increasing the batch size for GF-S, overall lifetimes of
473 batch sub-processes may increase sufficiently to allow better utilization of parallelization.
474 However, since GF-C results for larger batch sizes clearly show possible performance
475 degradation for more complex calculations (e.g. due to suboptimal load balancing), eight
476 was considered as a ‘safe’ default which improves overall performance for both simple and
477 complex calculation scenarios (Figure 5b).

478

479 Utilizing native multithreading for FFM, SIR (without native batch mode) and MF yields only
480 relatively small reductions in their execution times (Figure 5b). Under optimum conditions
481 (6-8 threads), the most significant drop was observed for SIR-C (~40%), followed by FFM-S,
482 FFM-C and MF-C (~20%). These results suggest that native multithreading only yields partial
483 parallelization, which primarily occurs with complex input conditions. Note that *SIRIUS*
484 supports different linear programming solvers (*Gurobi* [110], *CPLEX* [111] and the default
485 *GLPK* [112]), which may influence overall performance and parallelization [113].

486 Nevertheless, a comparison between these solvers did not reveal significant changes with
487 our experimental conditions (Additional file 2: Figure S5). Combining the multiprocessing
488 functionality with native multithreading under optimum conditions (i.e. 6 parallel
489 processes/threads) only reduces execution times for SIR-C (Figure 5b). As such, both
490 performance improvements and scalability of the multiprocessing implementation of
491 *patRoon* appear highly effective at this stage.

492

493 The native batch modes of MC and SIR allow calculations from multiple inputs within a
494 single execution. This reduces the total number of tool executions, which may (1) lower the
495 accumulated overhead associated with starting and finishing tool executions and (2) hamper
496 effective parallelization from multiprocessing, especially if executions are less than the
497 available CPU cores. The combination of multiprocessing (optimum conditions) and native
498 batch mode was benchmarked with increasing number of inputs per tool execution (i.e. the
499 native batch size; Additional file 2: Figure S6). For MC, execution times were largely
500 unaffected by the input batch size if multiprocessing was disabled, which indicates a low
501 execution overhead. Lowest execution times were observed when multiprocessing was
502 enabled with small batch sizes ($\leq 25\%$ of the total inputs), which indicates a lack of native
503 parallelization support. In contrast, SIR showed significantly lower overall execution times
504 with increasing batch sizes (up to $\sim 7000\%$ and $\sim 320\%$ for SIR-S and SIR-C, respectively),
505 while enabling multiprocessing did not reduce execution times for batch sizes > 1 . These
506 results show that (1) SIR has a relative large execution overhead, which impairs
507 multiprocessing performance gains, and (2) supports effective native parallelized batch
508 execution. Thus, SIR performs most optimal if all calculations are performed within a single

509 execution. Similar to previous SIR benchmarks, no significant differences were found across
510 different linear solvers (Additional file 2: Figure S7). The results demonstrate that
511 multiprocessing may improve efficiency for batch calculations with tools with low execution
512 overhead and/or lack of native parallelization. Nonetheless, the dramatic improvement in
513 SIR calculation times when using the native batch mode indicates that software authors
514 should generally consider implementing native threaded batch mode functionality if large
515 batch calculations are an expected use case.

516

517 Finally, the implemented optimization strategies were tested for a complete *patRoon* NTA
518 workflow consisting of typical data processing steps and using all previously tested tools.
519 The chosen input conditions roughly fell in between the aforementioned “simple” and
520 “complex” conditions (see code in Additional file 4). Note that optimization strategies were
521 unavailable for some steps (e.g. grouping of features and collection of MS peak lists), and
522 native batch mode was not used in order to demonstrate the usefulness of multiprocessing
523 for tools that do not support this (e.g. other tools than MC and SIR and those potentially
524 available in future versions of *patRoon*). Regardless, the benchmarks revealed a reduction in
525 total run times of ~50% (from ~200 to ~100 minutes; Figure 5c). Since execution times of
526 each step may vary significantly, the inclusion of different combinations of steps may
527 significantly influence overall execution times.

528

529 The use of multiprocessing for all tools (except SIR), the implemented batch mode strategies
530 for GF and the use of the native batch mode supported by SIR were set as default in
531 *patRoon* with the determined optimal parameters from the benchmarks results. However,

532 the user can still freely configure all these options to potentially apply further optimizations
533 or otherwise (partially) disable parallelization to conserve system resources acquired by
534 *patRoon*.

535
536 As a final note, it is important to realize that these benchmarks display execution times that
537 also involve preparing and processing results and include other overhead such as process
538 creation from *R*. For this reason, a direct comparison with standalone execution of
539 investigated tools was not possible. Nevertheless, the various optimization strategies
540 employed by *patRoon* minimize such overhead, and the added parallelization functionality
541 often provide a clear advantage in efficiency when using typical CLI tools in an *R* based NTA
542 workflow, especially considering the now widespread availability of computing systems with
543 increasing numbers of cores.

544 ***Demonstration: suspect screening***

545 The previous section investigated several parallelization strategies implemented in *patRoon*
546 for efficient data processing. A common method in environmental NTA studies to increase
547 data processing efficiency and reducing the data complexity is by merely screening for
548 chemicals of interest. This section demonstrates such a suspect screening workflow with
549 *patRoon*, consisting of (a) raw data pre-treatment, (b) extracting, grouping and suspect
550 screening of feature data, and finally (c) annotating features to confirm their identity. During
551 the workflow several rule-based filters are applied to improve data quality. The ‘suspects’ in
552 this demonstration are, in fact, a set of compounds spiked to influent samples (Additional
553 file 3: Table S2), hence, they were used for validation purposes of the workflow. After

554 completion of the suspect screening workflow, several methods are demonstrated to
555 inspect the resulting data.

556 **Suspect screening: workflow**

557 The code described here can easily be generated with the `newProject()` function, which
558 automatically generates a ready-to-use R script based on user input (section “Visualization,
559 reporting and graphical interface”).

560

561 First, the *patRoon* R package is loaded and a `data.frame` is generated with the file
562 information of the sample analyses and their replicate and blank assignments. Next, this
563 information is used to centroid and convert the raw analyses files to the open mzML file
564 format, a necessary step for further processing.

```
library(patRoon)

# Generate analysis file information for all files in a directory,
# assign replicate group names to all triplicates and specify which
# should be used for blank subtraction.
anaInfo <- generateAnalysisInfo("../data",
                                groups = c(rep("blank", 3),
                                           rep("influent-A", 3),
                                           rep("effluent-A", 3),
                                           rep("influent-B", 3),
                                           rep("effluent-B", 3)),
                                blanks = "blank")

convertMSFiles(anaInfo = anaInfo, from = "thermo", to = "mzML",
565                         algorithm = "pwiz", centroid = "vendor")
```

566 The next step involves finding features and grouping them across samples. This example
567 uses the *OpenMS* algorithms and sets several algorithm specific parameters that were
568 manually optimized for the employed analytical instrumentation to optimize the workflow
569 output. Other algorithms (e.g. *enviPick*, *XCMS*) are easily selected by changing the
570 `algorithm` function parameter.

```
571
features <- findFeatures(anaInfo, algorithm = "openms",
                         noiseThrInt = 4E3,
                         chromFWHM = 3, minFWHM = 1, maxFWHM = 30,
                         chromSNR = 5, mzPPM = 5)
fGroups <- groupFeatures(features, algorithm = "openms")
```

572 Several rule-based filters are then applied for general data clean-up, followed by the
573 removal of sample blanks from the feature dataset.

```
574
fGroups <- filter(fGroups,
                   # minimum absolute feature intensity
                   absMinIntensity = 1E5,
                   # must be present in all replicates
                   relMinReplicateAbundance = 1,
                   # max relative standard deviation replicate intensities
                   maxReplicateIntRSD = 0.75,
                   # minimum feature intensity above blank
                   blankThreshold = 5,
                   # remove blank analyses afterwards
                   removeBlanks = TRUE)
```

575 Next, features are screened with a given suspect list, which is a CSV file read into a
576 `data.frame` containing the name, SMILES and (optionally) retention time for each suspect
577 (see Additional file 5). While the list in this demonstration is rather small (18 compounds,
578 see SX), larger lists containing several thousands of compounds such as those available on
579 the NORMAN network Suspect List Exchange [114] can also be used. The screening results
580 are returned in a `data.frame`, where each row is a hit (a suspect may occur multiple times)
581 containing the linked feature group identifier and other information such as detected *m/z*
582 and retention time (deviations). Finally, this table is used to transform the original feature
583 groups object (`fGroups`) by removing any unassigned features and tagging remainders by
584 their suspect name.

```
585
suspects <- read.csv("suspects.csv")
scr <- screenSuspects(fGroups, suspects, mzWindow = 0.002,
                      rtWindow = 6, adduct = "[M+H]+")
fGroupsSusp <- groupFeaturesScreening(fGroups, scr)
```

586 In the final step of this workflow annotation is performed, which consists of (a) generation
587 of MS peak list data, (b) general clean-up to only retain significant MS/MS mass peaks,
588 automatic annotation of (c) formulae and (d) chemical compounds, and (e) combining both
589 annotation data to improve ranking of candidate compounds. As with previous workflow
590 steps, the desired algorithms (*mzR*, *GenForm* and *MetFrag* in this example) are set using the
591 `algorithm` function parameter. Similarly, the compound database used by *MetFrag* (here
592 *CompTox* via a local CSV file obtained from [115]) can easily be changed to other databases
593 such as *PubChem*, *ChemSpider* or another local file.

```
594 mslists <- generateMSPeakLists(fGroupsSusp, "mzr",
  precursorMzWindow = 0.5)
mslists <- filter(mslists, relMSMSIntThr = 0.02, topMSMSPeaks = 10)

formulas <- generateFormulas(fGroupsSusp, "genform", mslists,
  adduct = "[M+H]+",
  elements = "CHNOPSClBr")
# Configure location of CompTox CSV file
options(patRoon.path.MetFragCompTox =
  "C:/CompTox_17March2019_SelectMetaData.csv")
compounds <- generateCompounds(fGroupsSusp, mslists, "metfrag",
  adduct = "[M+H]+",
  database = "comptox")
compounds <- addFormulaScoring(compounds, formulas, updateScore = TRUE)
```

595 **Suspect screening: data inspection**

596 All data generated during the workflow (e.g. features, peak lists, annotations) can be
597 inspected by overloads of common *R* methods.

```

# intensities for each feature in first group
> fGroups[[1]]
[1] 210235.3 242051.9 254323.8 260419.1 205407.0 261099.1      0.0
0.0      0.0      0.0      0.0      0.0

# averaged MS/MS peak list for feature group of carbamazepine suspect
> mslists[["Carbamazepine"]]$MSMS
mz      intensity      precursor
1: 192.0804 284478.607      FALSE
2: 193.0880 69396.510      FALSE
3: 194.0960 1126534.943      FALSE
4: 237.1019 5406.667      TRUE

# compound annotation data for all features (subset shown for clarity)
> as.data.frame(compounds)[1:5, 1:5]
group      explainedPeaks score neutralMass   SMILES
1 n-Methylbenzotriazole-1 4 12.268046 133.064 NC1=NC2=CC=CC=C2N1
2 n-Methylbenzotriazole-1 5 9.546212 133.064 CC1=CC2=C(NN=N2)C=C1
3 n-Methylbenzotriazole-1 5 6.722034 133.064 NC1=CC=C2NN=CC2=C1
4 n-Methylbenzotriazole-1 5 6.715495 133.064 CC1=C2NN=NC2=CC=C1
5 n-Methylbenzotriazole-1 4 6.483770 133.064 CN1N=NC2=CC=CC=C12

```

598

599 Furthermore, all workflow data can easily be subset with e.g. the *R* subset operator ("["),
 600 for instance, to perform a (hypothetical) prioritization of features that are most intense in
 601 the effluent samples.

```

# obtain table with replicate averaged feature intensities
> intTab <- as.data.frame(fGroupsSusp, average = TRUE)
> head(intTab[, 1:5] # show first 5 rows/columns
group      ret      mz      influent-A effluent-A
1 n-Methylbenzotriazole-1 600.6524 134.0709 2021597.7      0.0
2 n-Methylbenzotriazole-2 607.5665 134.0709 2399435.6 192759.6
3      Barbital 137.3162 185.0918 145150.0      0.0
4      Benzotriazole 478.6665 120.0553 1494092.0 190069.0
5      Carbamazepine 797.5051 237.1018 2849756.3      0.0
6      Carbendazim 378.8226 192.0764 504191.7      0.0

# obtain group names from the 5 highest intense features in either
# of the effluents
> top1 <- intTab$group[order(intTab[["effluent-A"]]),
                           decreasing = TRUE)[1:5]
> top2 <- intTab$group[order(intTab[["effluent-B"]]),
                           decreasing = TRUE)[1:5]
> top <- union(top1, top2)
> top
[1] "Metformin"           "Terbutylazine"
[3] "Triphenylphosphine oxide" "Melamine-2"
[5] "n-Methylbenzotriazole-2" "Benzotriazole"
[7] "n-Methylbenzotriazole-1" "Propranolol"

# subset original object
> fGroupsSusp <- fGroupsSusp[, top]

```

602

603 Visualization of data generated during the workflow is performed by various plotting
604 functions (see Figure 6).

```
# plot unique features in influents
plot(fGroups[rGroups == c("influent-A", "influent-B")],
      colourBy = "rGroups", onlyUnique = TRUE)
# all EICs for a feature group
plotEIC(fGroupsSusp[, "Terbutylazine"], colourBy = "rGroup")
plotSpec(compounds, index = 1, groupName = "Benzotriazole",
         mslists)
plotUpSet(fGroupsSusp)
plotChord(fGroupsSusp, average = TRUE)
plotVenn(fGroupsSusp, which = c("influent-B", "effluent-B"))
```

605
606 The final step in a *patRoon* NTA workflow involves automatic generation of comprehensive
607 reports of various formats which allow (interactive) exploration of all data (see Additional
608 file 2: Figure S8).

```
reportCSV(fGroupsSusp, formulas = formulas, compounds = compounds)
reportPDF(fGroupsSusp, formulas = formulas,
          compounds = compounds, MSpeakLists = mslists)
reportHTML(fGroupsSusp, formulas = formulas,
           compounds = compounds, MSpeakLists = mslists)
```

609
610 **Suspect screening: results**
611 A summary of data generated during the NTA workflow demonstrated here is shown in
612 Table 5 and Additional file 3: Table S2. The complete workflow finished in approximately 8
613 minutes (employing a laptop with an Intel® Core™ I7-8550U CPU, 16 gigabyte RAM, NVME
614 SSD and the Windows 10 Pro operating system). While nearly 60 000 features were grouped
615 into nearly 20 000 feature groups, the majority (97%, 678 remaining) were filtered out
616 during the various pre-treatment filter steps. Regardless, most suspects were found (17/18
617 attributed to 19/20 individual chromatographic peaks), and the missing suspect (aniline)
618 could be detected when lowering the intensity threshold of the `filter()` function used to
619 post-filter feature groups in the workflow. The majority of suspects (17) were annotated
620 with the correct chemical compound as first candidate, the two n-methylbenzotriazole

621 isomer suspects were ranked as second or fourth. Results for formulae assignments were
622 similar, with the exception of dimethomorph, where the formula was ranked in only the top
623 twenty-five (the candidate chemical compound was ranked first, however).

624

625 While this demonstration conveys a relative simple NTA with 'known suspects', the results
626 show that *patRoon* (a) allows a straightforward approach to perform a complete and
627 tailored NTA workflow, (b) provides powerful general data clean-up functionality to
628 prioritize data and (c) realizes effective automated annotation of detected features.

629

630 *Table 5. Summarizing results for the demonstrated patRoon NTA workflow.*

		Amount
Features	Total found	57 113 (mean 3,808/sample)
Feature groups	Raw dataset	19 970
	Replicate filters (1 st pass ¹)	4 719 (-76%)
	Blank filter	2 933 (-85%)
	Intensity filters	964 (-95%)
	Replicate filters (2 nd pass ¹)	678 (-97%)
Suspects	Total found	19 out of 20
	Annotated	19
Formulae	Total candidates	163 (mean 9/feature group)
	Correctly ranked 1 st	13 (68%)
	Correctly ranked 1 st -2 nd	16 (84%)
	Correctly ranked 1 st -5 th	17 (89%)
Compounds	Total candidates	1 017 (mean 54/feature group)
	Correctly ranked 1 st	17 (85%)
	Correctly ranked 1 st -2 nd	18 (90%)
	Correctly ranked 1 st -5 th	19 (100%)

(1): Replicate filters are repeated if necessary, see section "Data reduction, comparison and conversion".

631 **Conclusions**

632 This paper presents *patRoon*, a fully open source platform that provides a comprehensive
633 MS based NTA data processing workflow developed in the *R* environment. Major workflow

634 functionality is implemented through the usage of existing and well-tested software tools,
635 connecting primarily open and a few closed approaches. The workflows are easily setup for
636 common use cases, while full customization and mixing of algorithms allows for execution of
637 completely tailored workflows. In addition, extensive functionality related to data
638 processing, annotation, visualization, reporting and others was implemented in *patRoon* to
639 provide an important toolbox for effectively handling complex NTA studies. The easy and
640 predictable interface of *patRoon* lowers the computational expertise required of users,
641 making it available for a broad audience. Major implemented optimization strategies were
642 demonstrated to reduce computational times. Furthermore, a typical suspect screening
643 workflow was demonstrated on real-world data from an environmental study related to
644 drinking water treatment.

645
646 *patRoon* has been under development for several years and has already been applied in a
647 variety of studies, such as the characterization of organic matter [71], elucidation of
648 transformation products of biocides [7, 12] and assessment of removal of polar organics
649 by reversed-osmosis drinking water treatment [14]. *patRoon* will undergo further
650 development, and extension of integrated workflow algorithms is planned for new and less
651 commonly used ones, while additional componentization strategies will be implemented to
652 help prioritizing data. Addition of new workflow functionality is foreseen, such as usage of
653 ion-mobility spectrometry data to assist annotation, automated screening of transformation
654 products (e.g. utilizing tools such as *BioTransformer* [116]), prediction of feature quantities
655 for prioritization purposes (recently reviewed in [117]) and automated chemical
656 classification (e.g. through *ClassyFire* [118]). Finally, interfacing with other *R* based mass

657 spectrometry software such as those provided by the “R for Mass Spectrometry” initiative
658 [119] is planned to further improve the interoperability of *patRoon*. The use in real-world
659 studies, feedback from users and developments within the non-target analysis community,
660 are all critical in determining future directions and improvements of *patRoon*. We envisage
661 that the open availability, straightforward usage, vendor independence and comprehensive
662 functionality will be useful to the community and result in a broad adoption of *patRoon*.

663 **Availability and requirements**

664 **Project name:** patRoon

665 **Project home page:** <https://github.com/rickhelmus/patRoon>

666 **Operating system(s):** Platform independent (tested on Microsoft Windows and Linux)

667 **Programming language(s):** R, C++, JavaScript

668 **Other requirements:** Depending on utilized algorithms (see installation instructions in [85,
669 88])

670 **License:** GNU GPL version 3

671 **Any restrictions to use by non-academics:** none

672 **Abbreviations**

673 **CEC:** Chemical of emerging concern

674 **CLI:** Command-line interface

675 **CMP:** Compound annotation

676 **CMT:** Componentization

677 **CSV:** Comma-separated value

678 **DBI:** The database interface

679 **EIC:** Extracted ion chromatogram

680 **FFM(-S/C):** FeatureFinderMetabo (simple/complex conditions)

681 **FG:** Feature groups

682 **FOR:** Formula annotation

683 **FTS:** Features

684 **GC:** Gas chromatography

685 **GC-MS:** GC coupled to mass spectrometry

686 **GF(-S/C):** GenForm (simple/complex conditions)

687 **HTML:** Hypertext markup language

688 **HRMS:** High resolution mass spectrometry

689 **IPO:** Isotopologue parameter optimization

690 **LC:** Liquid chromatography

691 **LC-MS:** LC coupled to mass spectrometry

692 **MC(-S/C):** msConvert (simple/complex conditions)

693 **MF(-S/C):** MetFrag (simple/complex conditions)

694 **MS/MS:** Tandem mass spectrometry

695 **MSPL:** MS peak list

696 **NTA:** Non-target analysis

697 **PDF:** Portable document format

698 **PO:** Parameter optimization

699 **PT:** Pre-treatment

700 **SIR(-S/C):** SIRIUS (simple/complex conditions)

701 **SUS:** Suspect screening
702 **XCMS:** Various forms (X) of chromatography mass spectrometry (*R* package MS data
703 processing)

704 **Definitions**

705 **Features (FTS):** data points assigned with unique chromatographic and mass spectral
706 information (e.g. retention time, peak area and accurate m/z), which potentially described a
707 compound in a sample analysis.

708 **Feature group (FG):** A group of features considered equivalent across sample analyses.

709 **MS peak list (MSPL):** tabular data (m/z and intensity) for MS or MS/MS peaks attributed to a
710 feature and used as input data for annotation purposes.

711 **Formula/Compound (FOR/CMP):** a chemical formula or compound candidate revealed
712 during feature annotation.

713 **Component (CMT):** A collection of feature groups that are somehow linked, such as MS
714 adducts, homologous series or highly similar intensity trends.

715 **Declarations**

716 ***Availability of data and materials***

717 The source code of *patRoon* and online versions of its manuals are available for download
718 from <https://github.com/rickhelmus/patRoon> and archived in [120]. The raw data used for
719 benchmarking and demonstration purposes in this manuscript is archived in [107]. The
720 scripts used to perform benchmarking and the input suspect list for demonstration purposes
721 are provided as Additional file 4 and 5, respectively.

722 ***Competing interests***

723 The authors declare that they have no competing interests.

724 ***Funding***

725 This work was internally funded by the Institute of Biodiversity and Ecosystem Dynamics
726 (University of Amsterdam). ELS is supported by the Luxembourg National Research Fund
727 (FNR) for project A18/BM/12341006.

728 ***Authors' contributions***

729 RH wrote the manuscript, source code, designed the experiments and interpreted the
730 results. ELS provided valuable feedback to improve the software. ELS and other authors
731 supervised this work and contributed to writing the manuscript. All authors read and
732 approved the final manuscript.

733 ***Acknowledgements***

734 The many authors involved in the open mass spectrometry software development
735 community are highly acknowledged as their contributions are the foundation for the
736 development of *patRoOn*. In addition, Vittorio Albergamo, Andrea Brunner, Thomas Wagner,
737 Olaf Brock and other users of *patRoOn* are thanked for testing and providing feedback for
738 future developments. We thank the Dutch drinking water companies Dunea and PWN for
739 sharing the raw HRMS data that was used for benchmarking and demonstration purposes.
740 Markus Fleischauer is acknowledged for his feedback on execution of batch execution of
741 *SIRIUS*. Finally, Olaf Brock is acknowledged for the design of some of the visualizations of
742 benchmarking data.

743 **References**

744 1. Hollender J, Schymanski EL, Singer HP, Ferguson PL (2017) Nontarget Screening with
745 High Resolution Mass Spectrometry in the Environment: Ready to Go? *Environ Sci
746 Technol* 51:11505–11512 . <https://doi.org/10.1021/acs.est.7b02184>

747 2. Chiaia-Hernandez AC, Schymanski EL, Kumar P, Singer HP, Hollender J (2014) Suspect
748 and nontarget screening approaches to identify organic contaminant records in lake
749 sediments. *Anal Bioanal Chem* 406:7323–7335 . [https://doi.org/10.1007/s00216-014-8166-0](https://doi.org/10.1007/s00216-014-
750 8166-0)

751 3. Sjerps RMA, Vughs D, van Leerdam JA, ter Laak TL, van Wezel AP (2016) Data-driven
752 prioritization of chemicals for various water types using suspect screening LC-HRMS.
753 *Water Research* 93:254–264 . <https://doi.org/10.1016/j.watres.2016.02.034>

754 4. Chiaia-Hernández AC, Günthardt BF, Frey MP, Hollender J (2017) Unravelling
755 Contaminants in the Anthropocene Using Statistical Analysis of Liquid
756 Chromatography–High-Resolution Mass Spectrometry Nontarget Screening Data
757 Recorded in Lake Sediments. *Environ Sci Technol* 51:12547–12556 .
758 <https://doi.org/10.1021/acs.est.7b03357>

759 5. Albergamo V, Schollée JE, Schymanski EL, Helmus R, Timmer H, Hollender J, de Voogt P
760 (2019) Nontarget Screening Reveals Time Trends of Polar Micropollutants in a
761 Riverbank Filtration System. *Environ Sci Technol* 53:7584–7594 .
762 <https://doi.org/10.1021/acs.est.9b01750>

763 6. Hernández F, Bakker J, Bijlsma L, de Boer J, Botero-Coy AM, Bruinen de Bruin Y, Fischer
764 S, Hollender J, Kasprzyk-Hordern B, Lamoree M, López FJ, Laak TL ter, van Leerdam JA,
765 Sancho JV, Schymanski EL, de Voogt P, Hogendoorn EA (2019) The role of analytical

766 chemistry in exposure science: Focus on the aquatic environment. *Chemosphere*
767 222:564–583 . <https://doi.org/10.1016/j.chemosphere.2019.01.118>

768 7. Wagner TV, Helmus R, Quiton Tapia S, Rijnaarts HHM, de Voogt P, Langenhoff AAM,
769 Parsons JR (2020) Non-target screening reveals the mechanisms responsible for the
770 antagonistic inhibiting effect of the biocides DBNPA and glutaraldehyde on benzoic
771 acid biodegradation. *Journal of Hazardous Materials* 386:121661 .
772 <https://doi.org/10.1016/j.jhazmat.2019.121661>

773 8. Kolkman A, Martijn BJ, Vughs D, Baken KA, van Wezel AP (2015) Tracing Nitrogenous
774 Disinfection Byproducts after Medium Pressure UV Water Treatment by Stable Isotope
775 Labeling and High Resolution Mass Spectrometry. *Environ Sci Technol* 49:4458–4465 .
776 <https://doi.org/10.1021/es506063h>

777 9. Schollée JE, Schymanski EL, Avak SE, Loos M, Hollender J (2015) Prioritizing Unknown
778 Transformation Products from Biologically-Treated Wastewater Using High-Resolution
779 Mass Spectrometry, Multivariate Statistics, and Metabolic Logic. *Anal Chem* 87:12121–
780 12129 . <https://doi.org/10.1021/acs.analchem.5b02905>

781 10. Brunner AM, Vughs D, Siegers W, Bertelkamp C, Hofman-Caris R, Kolkman A, ter Laak T
782 (2019) Monitoring transformation product formation in the drinking water treatments
783 rapid sand filtration and ozonation. *Chemosphere* 214:801–811 .
784 <https://doi.org/10.1016/j.chemosphere.2018.09.140>

785 11. Brunner AM, Bertelkamp C, Dingemans MML, Kolkman A, Wols B, Harmsen D, Siegers
786 W, Martijn BJ, Oorthuizen WA, ter Laak TL (2020) Integration of target analyses, non-
787 target screening and effect-based monitoring to assess OMP related water quality

788 changes in drinking water treatment. *Science of The Total Environment* 705:135779 .

789 <https://doi.org/10.1016/j.scitotenv.2019.135779>

790 12. Wagner TV, Helmus R, Becker E, Rijnaarts HHM, Voogt P de, Langenhoff AAM, Parsons

791 JR (2020) Impact of transformation, photodegradation and interaction with

792 glutaraldehyde on the acute toxicity of the biocide DBNPA in cooling tower water.

793 *Environ Sci: Water Res Technol* 6:1058–1068 . <https://doi.org/10.1039/C9EW01018A>

794 13. Jonker W, Lamoree MH, Houtman CJ, Hamers T, Somsen GW, Kool J (2015) Rapid

795 activity-directed screening of estrogens by parallel coupling of liquid chromatography

796 with a functional gene reporter assay and mass spectrometry. *Journal of*

797 *Chromatography A* 1406:165–174 . <https://doi.org/10.1016/j.chroma.2015.06.012>

798 14. Albergamo V, Escher BI, Schymanski EL, Helmus R, Dingemans MML, Cornelissen ER,

799 Kraak MHS, Hollender J, Voogt P de (2019) Evaluation of reverse osmosis drinking

800 water treatment of riverbank filtrate using bioanalytical tools and non-target

801 screening. *Environ Sci: Water Res Technol* 6:103–116 .

802 <https://doi.org/10.1039/C9EW00741E>

803 15. Brunner AM, Dingemans MML, Baken KA, van Wezel AP (2019) Prioritizing

804 anthropogenic chemicals in drinking water and sources through combined use of mass

805 spectrometry and ToxCast toxicity data. *Journal of Hazardous Materials* 364:332–338 .

806 <https://doi.org/10.1016/j.jhazmat.2018.10.044>

807 16. Zwart N, Jonker W, Broek R ten, de Boer J, Somsen G, Kool J, Hamers T, Houtman CJ,

808 Lamoree MH (2020) Identification of mutagenic and endocrine disrupting compounds

809 in surface water and wastewater treatment plant effluents using high-resolution

810 effect-directed analysis. Water Research 168:115204 .
811 <https://doi.org/10.1016/j.watres.2019.115204>

812 17. Schymanski EL, Singer HP, Slobodnik J, Ipolyi IM, Oswald P, Krauss M, Schulze T,
813 Haglund P, Letzel T, Grosse S, Thomaidis NS, Bletsou A, Zwiener C, Ibáñez M, Portolés
814 T, de Boer R, Reid MJ, Onghena M, Kunkel U, Schulz W, Guillon A, Noyon N, Leroy G,
815 Bados P, Bogialli S, Stipaničev D, Rostkowski P, Hollender J (2015) Non-target screening
816 with high-resolution mass spectrometry: critical review using a collaborative trial on
817 water analysis. *Anal Bioanal Chem* 407:6237–6255 . <https://doi.org/10.1007/s00216-015-8681-7>

819 18. Peisl BYL, Schymanski EL, Wilmes P (2018) Dark matter in host-microbiome
820 metabolomics: Tackling the unknowns—A review. *Analytica Chimica Acta* 1037:13–27 .
821 <https://doi.org/10.1016/j.aca.2017.12.034>

822 19. Martens L, Chambers M, Sturm M, Kessner D, Levander F, Shofstahl J, Tang WH,
823 Römpp A, Neumann S, Pizarro AD, Montecchi-Palazzi L, Tasman N, Coleman M,
824 Reisinger F, Souda P, Hermjakob H, Binz P-A, Deutsch EW (2011) mzML—a Community
825 Standard for Mass Spectrometry Data. *Molecular & Cellular Proteomics* 10: .
826 <https://doi.org/10.1074/mcp.R110.000133>

827 20. Pedrioli PGA, Eng JK, Hubley R, Vogelzang M, Deutsch EW, Raught B, Pratt B, Nilsson E,
828 Angeletti RH, Apweiler R, Cheung K, Costello CE, Hermjakob H, Huang S, Julian RK, Kapp
829 E, McComb ME, Oliver SG, Omenn G, Paton NW, Simpson R, Smith R, Taylor CF, Zhu W,
830 Aebersold R (2004) A common open representation of mass spectrometry data and its
831 application to proteomics research. *Nat Biotechnol* 22:1459–1466 .
832 <https://doi.org/10.1038/nbt1031>

833 21. Urban J, Afseth NK, Štys D (2014) Fundamental definitions and confusions in mass
834 spectrometry about mass assignment, centroiding and resolution. *TrAC Trends in*
835 *Analytical Chemistry* 53:126–136 . <https://doi.org/10.1016/j.trac.2013.07.010>

836 22. Chambers MC, Maclean B, Burke R, Amodei D, Ruderman DL, Neumann S, Gatto L,
837 Fischer B, Pratt B, Egertson J, Hoff K, Kessner D, Tasman N, Shulman N, Frewen B, Baker
838 TA, Brusniak M-Y, Paulse C, Creasy D, Flashner L, Kani K, Moulding C, Seymour SL,
839 Nuwaysir LM, Lefebvre B, Kuhlmann F, Roark J, Rainer P, Detlev S, Hemenway T,
840 Huhmer A, Langridge J, Connolly B, Chadick T, Holly K, Eckels J, Deutsch EW, Moritz RL,
841 Katz JE, Agus DB, MacCoss M, Tabb DL, Mallick P (2012) A cross-platform toolkit for
842 mass spectrometry and proteomics. *Nat Biotechnol* 30:918–920 .
843 <https://doi.org/10.1038/nbt.2377>

844 23. PubChem National Center for Biotechnology Information PubChem Database.
845 <https://pubchem.ncbi.nlm.nih.gov/>. Accessed 6 Feb 2020

846 24. Williams AJ, Grulke CM, Edwards J, McEachran AD, Mansouri K, Baker NC, Patlewicz G,
847 Shah I, Wambaugh JF, Judson RS, Richard AM (2017) The CompTox Chemistry
848 Dashboard: a community data resource for environmental chemistry. *Journal of*
849 *Cheminformatics* 9:61 . <https://doi.org/10.1186/s13321-017-0247-6>

850 25. Blaženović I, Kind T, Torbašinović H, Obrenović S, Mehta SS, Tsugawa H, Wermuth T,
851 Schauer N, Jahn M, Biedendieck R, Jahn D, Fiehn O (2017) Comprehensive comparison
852 of in silico MS/MS fragmentation tools of the CASMI contest: database boosting is
853 needed to achieve 93% accuracy. *Journal of Cheminformatics* 9:32 .
854 <https://doi.org/10.1186/s13321-017-0219-x>

855 26. Bruker MetaboScape. <https://www.bruker.com/products/mass-spectrometry-and-separations/ms-software/metaboscape.html>. Accessed 6 Feb 2020

856

857 27. Waters UNIFI Scientific Information System.

858 https://www.waters.com/waters/en_US/UNIFI-Scientific-Information-System/nav.htm?cid=134801359&locale=en_US. Accessed 6 Feb 2020

859

860 28. Thermo Scientific Compound Discoverer Software.

861 <https://www.thermofisher.com/uk/en/home/industrial/mass-spectrometry/liquid-chromatography-mass-spectrometry-lc-ms/lc-ms-software/multi-omics-data-analysis/compound-discoverer-software.html>. Accessed 6 Feb 2020

862

863

864 29. Progenesis QI. <http://www.nonlinear.com/progenesis/qi/>. Accessed 6 Feb 2020

865 30. Misra BB, Mohapatra S (2019) Tools and resources for metabolomics research

866 community: A 2017–2018 update. ELECTROPHORESIS 40:227–246 .

867 <https://doi.org/10.1002/elps.201800428>

868

869

870

871

872

873 31. Stanstrup J, Broeckling CD, Helmus R, Hoffmann N, Mathé E, Naake T, Nicolotti L, Peters K, Rainer J, Salek RM, Schulze T, Schymanski EL, Stravs MA, Thévenot EA, Treutler H, Weber RJM, Willighagen E, Witting M, Neumann S (2019) The

874 metaRbolomics Toolbox in Bioconductor and beyond. Metabolites 9:200 .

875 <https://doi.org/10.3390/metabo9100200>

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

999

876 33. Pluskal T, Castillo S, Villar-Briones A, Orešić M (2010) MZmine 2: Modular framework
877 for processing, visualizing, and analyzing mass spectrometry-based molecular profile
878 data. *BMC Bioinformatics* 11:395 . <https://doi.org/10.1186/1471-2105-11-395>

879 34. Hohrenk LL, Itzel F, Baetz N, Tuerk J, Vosough M, Schmidt TC (2020) Comparison of
880 Software Tools for Liquid Chromatography–High-Resolution Mass Spectrometry Data
881 Processing in Nontarget Screening of Environmental Samples. *Anal Chem* 92:1898–
882 1907 . <https://doi.org/10.1021/acs.analchem.9b04095>

883 35. R Core Team (2019) R: A Language and Environment for Statistical Computing. R
884 Foundation for Statistical Computing, Vienna, Austria

885 36. Lange E, Tautenhahn R, Neumann S, Gröpl C (2008) Critical assessment of alignment
886 procedures for LC-MS proteomics and metabolomics measurements. *BMC
887 Bioinformatics* 9:375 . <https://doi.org/10.1186/1471-2105-9-375>

888 37. Niu W, Knight E, Xia Q, McGarvey BD (2014) Comparative evaluation of eight software
889 programs for alignment of gas chromatography–mass spectrometry chromatograms in
890 metabolomics experiments. *Journal of Chromatography A* 1374:199–206 .
891 <https://doi.org/10.1016/j.chroma.2014.11.005>

892 38. Myers OD, Sumner SJ, Li S, Barnes S, Du X (2017) Detailed Investigation and
893 Comparison of the XCMS and MZmine 2 Chromatogram Construction and
894 Chromatographic Peak Detection Methods for Preprocessing Mass Spectrometry
895 Metabolomics Data. *Anal Chem* 89:8689–8695 .
896 <https://doi.org/10.1021/acs.analchem.7b01069>

897 39. Hao L, Wang J, Page D, Asthana S, Zetterberg H, Carlsson C, Okonkwo OC, Li L (2018)
898 Comparative Evaluation of MS-based Metabolomics Software and Its Application to

899 Preclinical Alzheimer's Disease. Scientific Reports 8:9291 .

900 <https://doi.org/10.1038/s41598-018-27031-x>

901 40. Myers OD, Sumner SJ, Li S, Barnes S, Du X (2017) One Step Forward for Reducing False

902 Positive and False Negative Compound Identifications from Mass Spectrometry

903 Metabolomics Data: New Algorithms for Constructing Extracted Ion Chromatograms

904 and Detecting Chromatographic Peaks. *Anal Chem* 89:8696–8703 .

905 <https://doi.org/10.1021/acs.analchem.7b00947>

906 41. Schymanski EL, Neumann S (2013) CASMI: And the Winner is . . . *Metabolites* 3:412–

907 439 . <https://doi.org/10.3390/metabo3020412>

908 42. Allen F, Pon A, Wilson M, Greiner R, Wishart D (2014) CFM-ID: a web server for

909 annotation, spectrum prediction and metabolite identification from tandem mass

910 spectra. *Nucleic Acids Res* 42:W94–W99 . <https://doi.org/10.1093/nar/gku436>

911 43. Allen F, Greiner R, Wishart D (2015) Competitive fragmentation modeling of ESI-

912 MS/MS spectra for putative metabolite identification. *Metabolomics* 11:98–110 .

913 <https://doi.org/10.1007/s11306-014-0676-4>

914 44. Loos M (2016) nontarget: Detecting Isotope, Adduct and Homologue Relations in LC-

915 MS Data. <https://CRAN.R-project.org/package=nontarget>

916 45. Loos M (2016) enviPick: Peak Picking for High Resolution Mass Spectrometry Data.

917 <https://CRAN.R-project.org/package=enviPick>. Accessed 2 Oct 2018

918 46. Loos M (2018) enviMass version 3.5 LC-HRMS trend detection workflow - R package.

919 <https://doi.org/10.5281/zenodo.1213098>

920 47. Meringer M, Reinker S, Zhang J, Muller A MS/MS data improves automated
921 determination of molecular formulas by mass spectrometry. *MATCH Commun Math*
922 *Comput Chem* 259–290

923 48. Ruttkies C, Schymanski EL, Wolf S, Hollender J, Neumann S (2016) MetFrag relaunched:
924 incorporating strategies beyond in silico fragmentation. *Journal of Cheminformatics* 8:3
925 . <https://doi.org/10.1186/s13321-016-0115-9>

926 49. FOR-IDENT LC. <https://water.for-ident.org/#!home>. Accessed 7 Feb 2020

927 50. Tsugawa H, Cajka T, Kind T, Ma Y, Higgins B, Ikeda K, Kanazawa M, VanderGheynst J,
928 Fiehn O, Arita M (2015) MS-DIAL: data-independent MS/MS deconvolution for
929 comprehensive metabolome analysis. *Nat Methods* 12:523–526 .
930 <https://doi.org/10.1038/nmeth.3393>

931 51. Tsugawa H, Kind T, Nakabayashi R, Yukihira D, Tanaka W, Cajka T, Saito K, Fiehn O,
932 Arita M (2016) Hydrogen Rearrangement Rules: Computational MS/MS Fragmentation
933 and Structure Elucidation Using MS-FINDER Software. *Anal Chem* 88:7946–7958 .
934 <https://doi.org/10.1021/acs.analchem.6b00770>

935 52. Röst HL, Sachsenberg T, Aiche S, Bielow C, Weisser H, Aicheler F, Andreotti S, Ehrlich H-
936 C, Gutenbrunner P, Kenar E, Liang X, Nahnsen S, Nilse L, Pfeuffer J, Rosenberger G,
937 Rurik M, Schmitt U, Veit J, Walzer M, Wojnar D, Wolski WE, Schilling O, Choudhary JS,
938 Malmström L, Aebersold R, Reinert K, Kohlbacher O (2016) OpenMS: a flexible open-
939 source software platform for mass spectrometry data analysis. *Nature Methods*
940 13:741–748 . <https://doi.org/10.1038/nmeth.3959>

941 53. Broeckling CD, Afsar FA, Neumann S, Ben-Hur A, Prenni JE (2014) RAMClust: A Novel
942 Feature Clustering Method Enables Spectral-Matching-Based Annotation for
943 Metabolomics Data. *Anal Chem* 86:6812–6817 . <https://doi.org/10.1021/ac501530d>

944 54. Böcker S, Letzel MC, Lipták Z, Pervukhin A (2009) SIRIUS: decomposing isotope patterns
945 for metabolite identification. *Bioinformatics* 25:218–224 .
946 <https://doi.org/10.1093/bioinformatics/btn603>

947 55. Dührkop K, Shen H, Meusel M, Rousu J, Böcker S (2015) Searching molecular structure
948 databases with tandem mass spectra using CSI:FingerID. *PNAS* 112:12580–12585 .
949 <https://doi.org/10.1073/pnas.1509788112>

950 56. Dührkop K, Böcker S (2015) Fragmentation Trees Reloaded. In: Przytycka TM (ed)
951 Research in Computational Molecular Biology. Springer International Publishing, pp
952 65–79

953 57. Böcker S, Dührkop K (2016) Fragmentation trees reloaded. *Journal of Cheminformatics*
954 8:5 . <https://doi.org/10.1186/s13321-016-0116-8>

955 58. Dührkop K, Fleischauer M, Ludwig M, Aksенов AA, Melnik AV, Meusel M, Dorrestein
956 PC, Rousu J, Böcker S (2019) SIRIUS 4: a rapid tool for turning tandem mass spectra into
957 metabolite structure information. *Nat Methods* 16:299–302 .
958 <https://doi.org/10.1038/s41592-019-0344-8>

959 59. Kuhl C, Tautenhahn R, Böttcher C, Larson TR, Neumann S (2012) CAMERA: An
960 Integrated Strategy for Compound Spectra Extraction and Annotation of Liquid
961 Chromatography/Mass Spectrometry Data Sets. *Anal Chem* 84:283–289 .
962 <https://doi.org/10.1021/ac202450g>

963 60. Tautenhahn R, Patti GJ, Rinehart D, Siuzdak G (2012) XCMS Online: A Web-Based
964 Platform to Process Untargeted Metabolomic Data. *Anal Chem* 84:5035–5039 .
965 <https://doi.org/10.1021/ac300698c>

966 61. Bruker DataAnalysis. <https://www.bruker.com/>. Accessed 20 Mar 2020

967 62. Libiseller G, Dvorzak M, Kleb U, Gander E, Eisenberg T, Madeo F, Neumann S,
968 Trausinger G, Sinner F, Pieber T, Magnes C (2015) IPO: a tool for automated
969 optimization of XCMS parameters. *BMC Bioinformatics* 16:118 .
970 <https://doi.org/10.1186/s12859-015-0562-8>

971 63. Eliasson M, Rännar S, Madsen R, Donten MA, Marsden-Edwards E, Moritz T, Shockcor
972 JP, Johansson E, Trygg J (2012) Strategy for Optimizing LC-MS Data Processing in
973 Metabolomics: A Design of Experiments Approach. *Anal Chem* 84:6869–6876 .
974 <https://doi.org/10.1021/ac301482k>

975 64. Loos M, Singer H (2017) Nontargeted homologue series extraction from hyphenated
976 high resolution mass spectrometry data. *J Cheminform* 9:12 .
977 <https://doi.org/10.1186/s13321-017-0197-z>

978 65. Schollée JE, Bourgin M, von Gunten U, McArdell CS, Hollender J (2018) Non-target
979 screening to trace ozonation transformation products in a wastewater treatment train
980 including different post-treatments. *Water Research* 142:267–278 .
981 <https://doi.org/10.1016/j.watres.2018.05.045>

982 66. Csardi G, Nepusz T (2006) The igraph software package for complex network research.
983 *InterJournal Complex Systems*:1695

984 67. Almende B.V., Thieurmel B, Robert T (2019) visNetwork: Network Visualization using
985 “vis.js” Library. <https://CRAN.R-project.org/package=visNetwork>

986 68. Kujawinski EB, Behn MD (2006) Automated Analysis of Electrospray Ionization Fourier
987 Transform Ion Cyclotron Resonance Mass Spectra of Natural Organic Matter. *Anal
988 Chem* 78:4363–4373 . <https://doi.org/10.1021/ac0600306>

989 69. Koch BP, Dittmar T (2006) From mass to structure: an aromaticity index for high-
990 resolution mass data of natural organic matter. *Rapid Communications in Mass
991 Spectrometry* 20:926–932 . <https://doi.org/10.1002/rcm.2386>

992 70. Koch BP, Dittmar T (2016) From mass to structure: an aromaticity index for high-
993 resolution mass data of natural organic matter. *Rapid Communications in Mass
994 Spectrometry* 30:250–250 . <https://doi.org/10.1002/rcm.7433>

995 71. Brock O, Helmus R, Kalbitz K, Jansen B Non-target screening of leaf litter-derived
996 dissolved organic matter using liquid chromatography coupled to high-resolution mass
997 spectrometry (LC-QTOF-MS). *European Journal of Soil Science.*
998 <https://doi.org/10.1111/ejss.12894>

999 72. Heller SR, McNaught A, Pletnev I, Stein S, Tchekhovskoi D (2015) InChI, the IUPAC
1000 International Chemical Identifier. *Journal of Cheminformatics* 7:23 .
1001 <https://doi.org/10.1186/s13321-015-0068-4>

1002 73. Guha R (2007) Chemical Informatics Functionality in R. *Journal of Statistical Software*
1003 18:1–16

1004 74. Schymanski EL, Gerlich M, Ruttkies C, Neumann S (2014) Solving CASMI 2013 with
1005 MetFrag, MetFusion and MOLGEN-MS/MS. *Mass Spectrometry* 3:S0036–S0036 .
1006 <https://doi.org/10.5702/massspectrometry.S0036>

1007 75. Langfelder P, Zhang B (2016) dynamicTreeCut: Methods for Detection of Clusters in
1008 Hierarchical Clustering Dendrograms. [https://CRAN.R-
1009 project.org/package=dynamicTreeCut](https://CRAN.R-project.org/package=dynamicTreeCut)

1010 76. Royal Society of Chemistry ChemSpider. <http://www.chemspider.com>. Accessed 6 Feb
1011 2020

1012 77. Lex A, Gehlenborg N, Strobelt H, Vuillemot R, Pfister H (2014) UpSet: Visualization of
1013 Intersecting Sets. *IEEE Transactions on Visualization and Computer Graphics* 20:1983–
1014 1992 . <https://doi.org/10.1109/TVCG.2014.2346248>

1015 78. Chen H, Boutros PC (2011) VennDiagram: a package for the generation of highly-
1016 customizable Venn and Euler diagrams in R. *BMC Bioinformatics* 12:35 .
1017 <https://doi.org/10.1186/1471-2105-12-35>

1018 79. Gu Z, Gu L, Eils R, Schlesner M, Brors B (2014) circlize implements and enhances
1019 circular visualization in R. *Bioinformatics* 30:2811–2812

1020 80. Gehlenborg N (2019) UpSetR: A More Scalable Alternative to Venn and Euler Diagrams
1021 for Visualizing Intersecting Sets. <https://CRAN.R-project.org/package=UpSetR>

1022 81. Xie Y, Allaire JJ, Grolemund G (2018) R Markdown: The Definitive Guide. Chapman and
1023 Hall/CRC, Boca Raton, Florida

1024 82. Allaire JJ, Xie Y, McPherson J, Luraschi J, Ushey K, Atkins A, Wickham H, Cheng J, Chang
1025 W, Iannone R (2019) rmarkdown: Dynamic Documents for R

1026 83. Iannone R, Allaire JJ, Borges B (2018) flexdashboard: R Markdown Format for Flexible
1027 Dashboards. <https://CRAN.R-project.org/package=flexdashboard>

1028 84. Chang W, Cheng J, Allaire JJ, Xie Y, McPherson J (2019) shiny: Web Application
1029 Framework for R. <https://CRAN.R-project.org/package=shiny>

1030 85. Helmus R (2020) patRoon 1.0.0 manuals. Zenodo.
1031 <https://doi.org/10.5281/zenodo.3889937>
1032 86. patRoon reference. <https://rickhelmus.github.io/patRoon/reference/index.html>. Accessed 11 Jun 2020
1034 87. patRoon tutorial. <https://rickhelmus.github.io/patRoon/articles/tutorial.html>. Accessed 11 Jun 2020
1036 88. Helmus R patRoon handbook.
1037 https://rickhelmus.github.io/patRoon/handbook_bd/index.html. Accessed 11 Jun 2020
1038 89. Xie Y (2016) bookdown: Authoring Books and Technical Documents with R Markdown.
1039 Chapman and Hall/CRC, Boca Raton, Florida
1040 90. Xie Y (2019) bookdown: Authoring Books and Technical Documents with R Markdown
1041 91. Wickham H, Danenberg P, Csárdi G, Eugster M (2019) roxygen2: In-Line Documentation
1042 for R. <https://CRAN.R-project.org/package=roxygen2>
1043 92. Helmus R (2020) patRoonData. <https://github.com/rickhelmus/patRoonData>. Accessed
1044 18 Mar 2020
1045 93. Helmus R, Albergamo V (2020) patRoonData: 1.0.0. Zenodo.
1046 <https://doi.org/10.5281/zenodo.3743266>
1047 94. Lang M (2017) checkmate: Fast Argument Checks for Defensive R Programming. The R
1048 Journal 9:437–445
1049 95. Csárdi G, Chang W (2019) processx: Execute and Control System Processes.
1050 <https://CRAN.R-project.org/package=processx>
1051 96. R Special Interest Group on Databases (R-SIG-DB), Wickham H, Müller K (2019) DBI: R
1052 Database Interface. <https://CRAN.R-project.org/package=DBI>

1053 97. Müller K, Wickham H, James DA, Falcon S (2019) RSQLite: “SQLite” Interface for R.
1054 <https://CRAN.R-project.org/package=RSQLite>

1055 98. Eddelbuettel D, François R (2011) Rcpp: Seamless R and C++ Integration. *Journal of*
1056 *Statistical Software* 40:1–18 . <https://doi.org/10.18637/jss.v040.i08>

1057 99. Eddelbuettel D (2013) Seamless R and C++ Integration with Rcpp. Springer, New York

1058 100. Eddelbuettel D, Balamuta JJ (2017) Extending R with C++: A Brief Introduction to Rcpp.
1059 *PeerJ Preprints* 5:e3188v1 . <https://doi.org/10.7287/peerj.preprints.3188v1>

1060 101. Kapoulkine A pugixml. <https://pugixml.org/>. Accessed 6 Feb 2020

1061 102. Dowle M, Srinivasan A (2019) data.table: Extension of `data.frame`. <https://CRAN.R-project.org/package=data.table>

1063 103. MetFragR. <http://ipb-halle.github.io/MetFrag/projects/metfragr/>. Accessed 6 Feb 2020

1064 104. Lang DT (2019) RDCOMClient: R-DCOM client

1065 105. Wickham H (2011) testthat: Get Started with Testing. *The R Journal* 3:5–10

1066 106. Henry L, Sutherland C, Hong D, Luciani TJ, Decorde M, Lise V (2019) vdiffr: Visual
1067 Regression Testing and Graphical Differing. <https://CRAN.R-project.org/package=vdiffr>

1068 107. Helmus R (2020) patRoon benchmarking & demonstration data. Zenodo.
1069 <https://doi.org/10.5281/zenodo.3885448>

1070 108. Mersmann O (2019) microbenchmark: Accurate Timing Functions. <https://CRAN.R-project.org/package=microbenchmark>

1072 109. Fischer B, Neumann S, Gatto L, Kou Q, Rainer J (2020) mzR: parser for netCDF, mzXML,
1073 mzData and mzML and mzIdentML files (mass spectrometry data).
1074 <https://bioconductor.org/packages/mzR/>. Accessed 6 Apr 2020

1075 110. Gurobi. <https://www.gurobi.com/>. Accessed 6 Feb 2020

1076 111. CPLEX Optimizer. <https://www.ibm.com/analytics/cplex-optimizer>. Accessed 6 Feb

1077 2020

1078 112. GNU Project - Free Software Foundation (FSF) GLPK (GNU Linear Programming Kit).

1079 <https://www.gnu.org/software/glpk/>. Accessed 6 Feb 2020

1080 113. Böcker S, Dührkop K, Fleischauer M, Ludwig M (2019) SIRIUS Documentation Release

1081 4.0.1

1082 114. NORMAN Suspect List Exchange – NORMAN SLE. [https://www.norman-](https://www.norman-network.com/nds/SLE/)

1083 [network.com/nds/SLE/](https://www.norman-network.com/nds/SLE/). Accessed 13 Mar 2020

1084 115. CompTox March 2019 CSV file.

1085 ftp://newftp.epa.gov/COMPTOX/Sustainable_Chemistry_Data/Chemistry_Dashboard/

1086 MetFrag_metadata_files/CompTox_17March2019_SelectMetaData.csv

1087 116. Djoumbou-Feunang Y, Fiamoncini J, Gil-de-la-Fuente A, Greiner R, Manach C, Wishart

1088 DS (2019) BioTransformer: a comprehensive computational tool for small molecule

1089 metabolism prediction and metabolite identification. *Journal of Cheminformatics* 11:2 .

1090 <https://doi.org/10.1186/s13321-018-0324-5>

1091 117. Kruve A (2019) Semi-quantitative non-target analysis of water with liquid

1092 chromatography/high-resolution mass spectrometry: How far are we? *Rapid*

1093 *Communications in Mass Spectrometry* 33:54–63 . <https://doi.org/10.1002/rcm.8208>

1094 118. Djoumbou Feunang Y, Eisner R, Knox C, Chepelev L, Hastings J, Owen G, Fahy E,

1095 Steinbeck C, Subramanian S, Bolton E, Greiner R, Wishart DS (2016) ClassyFire:

1096 automated chemical classification with a comprehensive, computable taxonomy. *J*

1097 *Cheminform* 8:61 . <https://doi.org/10.1186/s13321-016-0174-y>

1098 119. R for Mass Spectrometry. www.rformassspectrometry.org. Accessed 13 Mar 2020

1099 120. Rick Helmus (2020) patRoon 1.0.0. Zenodo. <https://doi.org/10.5281/zenodo.3889856>

1100 **Figures**

1101 **Figure 1. Generic workflow for environmental non-target analysis.**

1102

1103 **Figure 2. Overview of the NTA *patRoon* workflow.** All steps are optional. Steps that are
1104 connected by blue and straight arrows represent a one-way data dependency, whereas
1105 steps connected with red curved and dashed arrows represent steps with two-way data
1106 interaction.

1107

1108 **Figure 3. Graphical user interface tools in *patRoon*.** Tools are provided (a) to create a new
1109 *patRoon* data analysis project and (b) to inspect feature chromatography data.

1110

1111 **Figure 4. Interface for the *patRoon* workflow.** The workflow steps are performed by a set of
1112 functions that execute the selected algorithm and return the data in a harmonized format
1113 by utilizing the 'S4' object oriented programming approach of *R*. These objects all derive
1114 from a common base class and may be further sub-classed in algorithm specific classes (as is
1115 exemplified for features). Generic functions are defined for all workflow classes to
1116 implement further data processing functionality in a predictable and algorithm independent
1117 manner (see also Table 3). Further information is provided in the reference manual [85, 86].

1118

1119 **Figure 5. Parallelization benchmark results.** (a) Benchmark results for commonly used CLI
1120 tools applied in *patRoon* workflows under varying parallelization conditions. Tests were

1121 performed with “simple” (left) and “complex” (right) input conditions (Table 4) to simulate
1122 varying workflow complexity. Parallelization was performed with the multiprocessing
1123 functionality of *patRoon* (top) or by using native multithreading (bottom, for tools that
1124 supported this). Graphs represent number of processes or threads versus relative execution
1125 time (normalized to sequential results). The dotted grey lines represent the theoretical
1126 trend if maximum parallelization performance is achieved. The dashed blue line represents
1127 the number of physical cores that became the default selection in *patRoon* based on these
1128 results. (b) Comparison of execution times (normalized to the execution times of the
1129 unoptimized results) when tools are executed without optimizations (green), executed with
1130 native multithreading (FFM, SIR and MF) or batch mode (GF) (orange), executed with
1131 multiprocessing (purple) or a combination of the latter two (pink), using simple (left) and
1132 complex (right) input conditions. (c) Overview of execution times for a complete *patRoon*
1133 workflow executed under optimized versus unoptimized conditions. All results for MC and
1134 SIR were obtained without enabling their native batch mode.

1135

1136 **Figure 6. Common visualization functionality of *patRoon* applied to the demonstrated**
1137 **workflow.** From left to right: an *m/z* vs retention time plot of all feature groups, an EIC for
1138 the tramadol suspect found in both influent samples, a compound annotated spectrum for
1139 the 1,2,3-benzotriazole suspect and comparison of feature presence between sample
1140 groups using UpSet [77], Venn and chord diagrams.

1141 **Supplementary information**

1142 **Additional file 1:** Comma-separated file (.csv). Overview of software and databases that are
1143 used in the implementation in *patRoon*. This table summarizes all the software and
1144 databases that are described in the implementation section of the main text.

1145 **Additional file 2:** Word document (.docx). Supplementary figures. Additional figures that
1146 illustrate implementation details of *patRoon* and miscellaneous benchmarking results.

1147 **Additional file 3:** Word document (.docx). Supplementary tables. Additional tables with
1148 more details on the implementation and suspect screening demonstration.

1149 **Additional file 4:** Zip archive (.zip). Source code for benchmarks. Archive with several *R*
1150 scripts that were used to perform the parallelization benchmarks.

1151 **Additional file 5:** Comma-separated file (.csv). Demonstration suspect list. Suspect list that
1152 was used for the *patRoon* demonstration. The list was based on the detected compounds
1153 reported in [11], and SMILES identifiers for each suspect were collected from PubChem [23].

1154

	HRMS data	Features				Annotation						Primary Interface			License	References			
		PP	FTS	FG	C	SUS	MS	FA	CA	LA	HS	GA	C	RT					
a <i>CFM-ID</i>									X	X					CLI, Web C++	Cross	GPLv2.1	[42, 43]	
b <i>enviMass</i> , <i>enviPick</i> , <i>nontarget</i>	X ⁱ	X	X	X	X					X	X				GUI, R, Web	R	Cross	GPLv3.0 ¹	[44–46]
c <i>GenForm</i>							X								CLI	C++	Cross ²	GPLv2.0	[47]
d <i>MetFrag</i>								X	X			X	X		CLI, R, Web	Java	Cross	GPLv2.0	[48]
e <i>FOR-IDENT</i>								X ^d	X			X			Web	HTML	Cross	Closed	[49]
f <i>MS-DIAL</i> , <i>MS-FINDER</i>		X	X	X	X		X	X	X	X		X			CLI, GUI	C#	Win	GPLv3.0	[50, 51]
g <i>MZmine</i>	X	X ^{gl}	X	X	X		X	X ^k	X		X ^{gl}				GUI	Java	Cross	GPLv2.0	[33]
h <i>OpenMS</i>	X ^{hi}	X	X				X		X ^k	X		X			CLI, GUI, Python	C++	Win, Lin, Mac	BSD/3-Clause	[52]
i <i>ProteoWizard</i>	X														CLI, GUI	C++	Win, Lin	Apache 2.0	[22]
j <i>RAMClustR</i>							X				X				R	R	Cross	GPLv2.0	[53]
k <i>SIRIUS</i> and <i>CSI:FingerID</i>								X	X			X			CLI, GUI	Java	Cross	GPLv3.0	[54–58]
l <i>XCMS</i> and <i>CAMERA</i>		X	X	X							X				R	R	Cross	GPLv2.0	[32, 59]
m <i>XCMS Online</i>	X	X ^l	X ^l				X		X		X				Web	R	Cross	Closed	[60]
n <i>patRoon</i>	X ^{hi}	X ^{bhl}	X ^{hl}	X	X		X	X ^{ck}	X ^{dk}	X ^d	X ^b	X ^{il}	X	X ^d	R	R	Cross	GPLv3.0	

P: pre-processing; FTS: find features; FG: group features across samples; C: data clean up; SUS: suspect screening; MS: automatic MS data extraction for annotation purposes; FA: formula annotation; CA/LA: compound annotation (*in silico*/library); HS: unsupervised homologous series extraction; GA: grouping and annotating chemically related features (e.g. adducts, isotopes, in-source fragments); RT: retention time prediction; Bold: functionality integrated in *patRoon*; superscript: implemented with algorithms by given rows (omitted if only native); CLI: command-line interface; GUI: graphical user interface; Web: interfaced via internet browser; OS: Supported Operating Systems; Win: *Microsoft Windows*; (Lin): *GNU/Linux*, (Mac): *macOS*; Cross: cross-platform; (1): *enviMass* is distributed commercially; (2): Only *Microsoft Windows* binaries are distributed.