

Facilitating high resolution mass spectrometry data processing for screening of environmental water samples: An evaluation of two deconvolution tools

Richard Bade^{a,1}, Ana Causanilles^b, Erik Emke^b, Lubertus Bijlsma^a, Juan V. Sancho^a, Felix Hernandez^a, Pim de Voogt^{b,c*}

^a *Research Institute for Pesticides and Water, University Jaume I, Avda. Sos Baynat s/n, E-12071 Castellón, Spain.*

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

^c *Institute for Biodiversity and Ecosystem Dynamics, University of Amsterdam, P.O. Box 94248, 1090 GE, Amsterdam, The Netherlands*

* Corresponding author: w.p.devoogt@uva.nl Tel.: +31 20 5256565; fax: +31 20 5257431

¹ *Visiting researcher at KWR Watercycle Research Institute*

1 Abstract

2

3 A screening approach was applied to influent and effluent wastewater samples. After
4 injection in a LC-LTQ-Orbitrap, data analysis was performed using two deconvolution
5 tools, MsXerator (modules MPeaks and MS Compare) and Sieve 2.1. The outputs were
6 searched incorporating an in-house database of more than 200 pharmaceuticals and illicit
7 drugs or ChemSpider. This hidden target screening approach led to the detection of
8 numerous compounds including the illicit drug cocaine and its metabolite
9 benzoylecgonine and the pharmaceuticals carbamazepine, gemfibrozil and losartan. The
10 compounds found using both approaches were combined, and isotopic pattern and
11 retention time prediction were used to filter out false positives. The remaining potential
12 positives were reanalysed in MS/MS mode and their product ions were compared with
13 literature and/or mass spectral libraries. The inclusion of the chemical database
14 ChemSpider led to the tentative identification of several metabolites, including
15 paraxanthine, theobromine, theophylline and carboxylosartan, as well as the
16 pharmaceutical phenazone. The first three of these compounds are isomers and they were
17 subsequently distinguished based on their product ions and predicted retention times. This
18 work has shown that the use deconvolution tools facilitates non-target screening and
19 enables the identification of a higher number of compounds.

20

21 Keywords: Non-target screening, peak-picking, hidden target screening, software, high
22 resolution mass spectrometry, aquatic environment

23 *Introduction*

24 The investigation of emerging contaminants has become prevalent in analytical
25 environmental chemistry circles. The use of pharmaceuticals, personal care products and
26 illicit drugs is increasing worldwide, due to the growing population and the rise in
27 available products and the amount of these contaminants entering the aquatic
28 environment is of concern (Fatta-Kassinos et al., 2011). There is no blanket removal
29 process able to be undertaken by wastewater treatment plants (WWTPs) for all
30 compounds, leading to poor removal rates and detection of many of these compounds in
31 effluent wastewaters (EWW) and consequently in surface waters (Bijlsma et al., 2012;
32 Luo et al., 2014; van der Aa et al., 2013).

33 High Resolution Mass Spectrometry (HRMS) instruments, such as Time of Flight (TOF)
34 and Orbitrap have revolutionized the investigation of emerging contaminants in the
35 aquatic environment due to their high sensitivity in full scan mode, their increased mass
36 accuracy and the possibility to distinguish the isotopic pattern. HRMS instruments have
37 the ability to screen for unknowns due to exact mass measurements and these are unique
38 characteristics compared to other mass spectrometry instruments. Hybrid systems i.e.
39 HRMS hyphenated to a quadrupole or linear ion trap (LTQ), such as the LTQ-Orbitrap,
40 combines the tandem mass spectrometric capability associated with the LTQ with the
41 high mass resolving power (up to 100,000 FWHM) and mass accuracy capability of the
42 Orbitrap (de Voogt et al., 2011; Makarov and Scigelova, 2010). These hybrid
43 configurations based on HRMS allow reliable interpretation of MS/MS spectra and are
44 very valuable when dealing with complex environmental matrices, such as wastewater,
45 where co-elution of analytes with matrix interferences can result in ambiguous peaks
46 (Hogenboom et al., 2009). By utilising the ultra-high resolution capabilities, isobaric
47 compounds can easily be differentiated (Hernández et al., 2012).

48 In the literature, three different approaches are described for the detection and/or
49 identification of compounds: target, suspect/post-target and non-target (Aceña et al.,
50 2015; Bletsou et al., 2015; Gago-Ferrero et al., 2015; Hernández et al., 2015b, 2014, 2005;
51 Krauss et al., 2010; Leendert et al., 2015). Target methods are limited to a restricted
52 number of compounds, for all of which reference standards must be obtained and,
53 therefore, information on the occurrence of other unknown, relevant micropollutants may
54 be missed. Suspect screening takes advantage of a database of “known” compounds,
55 including molecular formulae, fragmentation and retention time, which can then be

56 computationally correlated to spectral HRMS data to give potential positive compounds.
57 As the concept of suspect screening implies that reference standards are not necessarily
58 available, the tentative identification of potential positives needs to be confirmed by the
59 use of reference standards (and MS/MS injections, if required) in a final step.

60 The third, non-target, approach is of increasing interest but notoriously difficult to
61 undertake, as, strictly speaking, no *a priori* information is available (Krauss et al., 2010;
62 Schymanski et al., 2014b; Zedda and Zwiener, 2012). Even with the help of automated
63 peak-picking software, thousands of peaks can be detected in an individual sample (Hug
64 et al., 2014; Kern et al., 2009). Consequently, subsequent steps must then be made to
65 reduce the number of peaks to a more manageable number, including molecular formula
66 derivation, isotopic pattern, mass defect analysis and retention time prediction (Gago-
67 Ferrero et al., 2015; Helbling et al., 2010; Kind and Fiehn, 2007). Further confidence in
68 the “potential positives” remaining can be gained through the use of fragmentation in a
69 subsequent MS/MS injection and comparison with *in silico* fragmentation and/or mass
70 spectral libraries (Bletsou et al., 2015; Gerlich and Neumann, 2013; Herrera-Lopez et al.,
71 2014; Hug et al., 2014; Little et al., 2012), with the latter referred to as “hidden targets”
72 (Letzel et al., 2015). In these situations, it is of prime importance and for ease of the
73 analyst to have software capable of fulfilling most (if not all) of these steps automatically.
74 Most manufacturers have software specific for their instrument and data, which can
75 automatically extract analytes of interest from the raw data, to facilitate suspect screening
76 approaches. However, despite the tremendous advances in software for
77 metabolite/transformation product detection and further non-target work, sometimes not
78 all required information is available in one platform, leading users to manufacturer-
79 independent software, such as the Eawag open-source R-code packages *enviMass*,
80 *enviPick*, *nontarget* and *RMassBank* (Schollée et al., 2015; Schymanski et al., 2014a)
81 which can enable the incorporation of additional parameters, such as the steps outlined
82 above. In spite of these problems, non-target screening is necessary to identify new or
83 unknown relevant pollutants, which is why efforts need to be made in developing proper
84 software and efficient identification tools.

85 This work portrays the combination of non-target data processing and hidden target
86 searching of environmental water samples after injection in an LC-LTQ-Orbitrap. Two
87 computational programs were utilized: *MsXelerator* (*MsMetrix*) and *Sieve 2.1* (Thermo
88 Scientific). An in-house database of more than 200 pharmaceuticals, personal care

89 products and illicit drugs was incorporated in both programs. Additionally, Sieve 2.1 was
90 used in combination with the ChemSpider search feature. The main objective was to
91 demonstrate the utility and additional value of these software packages for screening.
92 This led to the detection of numerous compounds across both programs. The compounds
93 detected by both methods were then reinjected to obtain MS/MS fragmentation, leading
94 to the tentative identification of 24 compounds. Ultimately, this work shows that the
95 combined use of two deconvolution tools combined with two hidden target screening
96 approaches provides more information than either one used individually.

97 *2. Materials and Methods*

98 *2.1 Reagents*

99 HPLC-grade methanol (MeOH), and formic acid (>98 % w/w) were purchased from
100 Mallinckrodt Baker (Deventer, The Netherlands). The ultrapure water was obtained by
101 purifying demineralized water in a Milli-Q system from Millipore (Bedford, MA, USA).
102 SPE cartridges used were Oasis HLB 3 mL (60 mg) from Waters (Milford, MA, USA).
103 Polytyrosine-1,3,6 standard used for mass axis calibration was purchased from Cs Bio
104 Co. (Menlo Park, CA, USA). Mixed cellulose ester membrane filters (0.45 µm) were
105 purchased from Whatman (Dassel, Germany).

106

107 *2.2 Water samples and extraction procedure*

108 Seven influent wastewater (IWW) and seven effluent wastewater (EWW) 24-hour
109 composite samples were collected over seven consecutive days in March 2014. They were
110 stored in high density polystyrene bottles, immediately centrifuged and stored in the dark
111 at -20°C. Analyses were performed as soon as possible after collection in order to keep
112 biotic or abiotic degradation to a minimum (Llorca et al. 2014).

113 A solid phase extraction (SPE) step was applied prior to analysis to pre-concentrate the
114 samples. All samples were filtered through a mixed cellulose ester membrane filter (0.45
115 µm). SPE was performed using Oasis HLB cartridges (60 mg). The water samples (EWW
116 100 mL, with IWW four times diluted (i.e. 25 mL sample diluted to 100mL by adding
117 Milli-Q water)) were loaded onto the cartridges and reconstructed in 1 mL of 10:90
118 MeOH:H₂O after elution with MeOH (5 mL). A procedural blank was also made,
119 following the steps above but using Milli-Q water. Analyses were performed by injecting
120 20 µL of the final extract (in triplicate) into the LC-LTQ FT Orbitrap. For further
121 information on the SPE procedure, see (Hernández et al., 2015a).

122 *2.3 Liquid Chromatography*

123 The HPLC system, consisted of a Surveyor auto sampler model Plus and a Surveyor
124 quaternary gradient HPLC-pump (Thermo Fisher Scientific, Breda, The Netherlands).
125 Chromatographic separation of the compounds was made using an XBridge C18 column
126 (150 mm × 2.1 mm I.D., particle size 3.5 µm) (Waters). The pre-column used was a 4.0
127 mm × 2.0 mm I.D. Phenomenex Security Guard column (Bester, Amsterdam, the

128 Netherlands). The analytical column and the guard column were maintained at a
129 temperature of 21°C in a column thermostat. A gradient was used at a constant flow rate
130 of 0.3 mL min⁻¹ using Milli-Q water (Solvent A) and MeOH (Solvent B) both with 0.05%
131 formic acid. The percentage of organic modifier (B) was changed linearly as follows: 0
132 min, 5%; 40 min, 100%; 45 min, 100%; 47 min, 5%. Between consecutive runs, the
133 analytical column was re-equilibrated for 5 min.

134 *2.4 LTQ-FT Orbitrap mass spectrometry*

135
136 An LTQ FT Orbitrap mass spectrometer (Thermo Electron, Bremen, Germany) was used.
137 The LTQ part of this system was equipped with a Heated Ion Max Electrospray Ionization
138 (HESI) probe and operated in the positive ion mode. The conditions were: source voltage
139 3.0 kV, heated capillary temperature 300°C, vaporizer temperature 350°C, capillary
140 voltage 13 V and tube lens 70 V. Products ions were generated in the LTQ trap at a
141 normalized collision energy setting of 35% and using an isolation width of 2 Da.

142

143 Full-scan accurate mass spectra (mass range from 50 to 1300 Da) were obtained at a mass
144 resolution of 60,000. The total cycle time depends upon the resolution; at the selected
145 resolution the total cycle time is 0.5 s. The instrument was initially set to operate in full-
146 scan ('survey') mode with accurate mass measurements. When an ion exceeded a preset
147 threshold, the instrument switched to product-ion scan mode in the ion trap part. Further
148 details on instrument operating conditions can be found elsewhere (Bijlsma et al., 2013)

149

150 All data were acquired and processed using Xcalibur version 2.1 software. A second
151 MS/MS injection was made by incorporating an inclusion list of masses (see Supporting
152 Information (S.I.) **Table S1** for list) with a retention window of ±2 minutes and collision
153 energy of 35%. Since MS/MS fragmentation was carried out in the ion trap, only nominal
154 mass was measured.

155

156 Mass axis calibration was performed with every batch run just prior to starting the batch
157 by using flow injection of a polytyrosine-1,3,6 solution ([M+H]⁺ 182.01170/508.20783
158 and 997.39781) at a flow rate of 10 µL min⁻¹.

159

160

161 *2.5 Settings of the Deconvolution Tools*

162 MsXelerator (MsMetrix)

163 MS Compare and MPeaks are modules within MsXelerator. MS Compare is specifically
164 designed for comparing MS spectra, whereas the MPeaks module picks peaks, with the
165 “keep largest C13 peaks only” and the “peak cluster” algorithms used to help discard
166 some, the latter performing componentization which groups together all peaks (i.e.
167 isotopes and adducts) arising from a single retention time (Table S2). All samples were
168 uploaded individually and later investigated as triplicates, corresponding to the three
169 triplicate injections of each sample. Procedural blank samples were initially processed
170 using the optimized software settings (see below) to subtract identical peaks from each
171 wastewater sample.

172 The “peak picking” was carried out by MPeaks on each individual sample using the
173 following parameters and values: Base peak width =11 (arbitrary units); spike width = 5
174 scans; peak separation = 5 scans; peak threshold = 0.5% (vs. largest peak); smoothness
175 threshold = 0.65%, signal/noise ratio = 20. The sensitivity value, which helps the user
176 find more or less sensitive parameters for the previous three parameters, was set relatively
177 low, at 2 (out of a maximum setting of 6). The peaks picked using these parameters were
178 further reduced by only keeping peaks relating to an $[M+H]^+$ charge state.

179 Using the second module, MS Compare, all samples were subjected to the following
180 LC/MS settings (in accurate mass mode) of the module for peak picking across multiple
181 samples: No baseline correction; FWHM (scans) = 3 (min) - 40 (max); min peak height
182 = 10,000 counts; delete spikes; m/z range: 100-650; Max. shift between peaks for
183 grouping: 20 scans; Time window for XIC: 0.25min; Mass accuracy: 10 ppm.

184 Sieve 2.1 (Thermo Scientific)

185 Sieve 2.1 combines the power of the two modules from MsXelerator. After an initial
186 peak-picking process using the settings described below, it compares MS spectra of the
187 procedural blank samples and the studied wastewater samples. Only compounds with an
188 m/z between 150 and 500, and only protonated molecules ($[M+H]^+$) were considered. The
189 list of potential positives were then search by either the in-house database or ChemSpider,
190 which were incorporated in the software. The results of the ChemSpider search were
191 exported into Microsoft Excel, and positive “hits” were considered based on their mass

192 error (< 2 ppm), and if the compound commercially existed i.e. hits which only
193 represented chemical formula were excluded.

194 The “control compare trend” feature of Sieve 2.1 (Thermo Scientific) was used with the
195 following parameters: peak intensity threshold = 250,000 (62,500 for IWW); m/z range
196 = 100-650; m/z width = 10ppm; retention time range = 3-40 minutes; maximum number
197 of frames = 5,000; frame time width = 1.00 minute; align bypass = true. For the hidden
198 target screening, either the database used in the analyses of MPeaks and MS Compare or
199 ChemSpider was incorporated.

200 *2.6 General workflow*

201 The general workflow followed in this work (pictorialized in Figure 1) falls into the
202 “hidden target” area of non-target screening, hypothesized by Letzel *et al.* (Letzel et al.,
203 2015), wherein non-target techniques (i.e. peak picking) are originally applied, but a
204 database (i.e. in-house database or ChemSpider) is used for identification.

205 All samples were injected in triplicate and the data were processed with the different
206 software packages of MsXelerator or Sieve 2.1. Only peaks in all three injections were
207 further investigated. This resulted in a list of chromatographic peaks, based on the
208 accurate masses of their protonated molecules. To gain a list of potential positives, two
209 hidden target identification methods were used: 1) an in-house database, containing more
210 than 200 parent compounds and metabolites and the online database ChemSpider. False
211 positives were manually removed after investigating the isotopic pattern (for the
212 characteristic patterns of sulfur- and chlorine-containing species) and retention time
213 prediction. A final “target list” was investigated by reinjecting the samples in MS/MS
214 mode, to get product ions. Fragmentation was then compared with online databases and
215 literature, which allowed the tentative identification of several compounds.

216

217 3. Results and Discussion

218 In this study, in order to show the progression through confidence levels of identification,
219 the terminology proposed in the literature by Hernández *et al.* (Hernández et al., 2015)
220 and Schymanski *et al.* (Schymanski et al., 2014) were followed. It must be noted that
221 potential positives and detected compounds, differentiated in this work, would both be
222 level 3 tentative candidate in the terminology of Schymanski *et al.* The final, tentatively
223 identified compounds are of a higher confidence level (level 2a). However, in order to
224 have total confirmation (level 1), reference standards are necessary. As no reference
225 standards were utilized in this work, this level could not be attained.

226 3.1 Optimization of the workflow

227 All samples were injected and processed in triplicate, which were compared together,
228 with only peaks in all three injections being further investigated. Procedural blank
229 samples were processed first, to subtract identical peaks from each subsequent IWW and
230 EWW sample.

231 The m/z range of MS Compare was made quite narrow as the compounds of interest in
232 this study and in the in-house database (small pharmaceutical/drug molecules) would be
233 within that range. The retention time range was reduced just to 3-40 mins to reduce the
234 likelihood of erroneously detecting species that elute very early and late due to the
235 high/low ratio of organic modifier, with the vast majority of all peaks in the total ion
236 chromatogram falling within this range. In spite of the known mass accuracy capability
237 of the Orbitrap, the mass accuracy was set at 10ppm to ensure that no compound would
238 be missed. After this processing, a list of masses common within each triplicate set was
239 made, with compounds being detected using the same peak peaking parameters and
240 database used in the final step of the MPeaks analysis.

241 Sieve 2.1 used the same m/z range, m/z width and retention time range as MS Compare
242 for better ease of results comparison. The peak intensity threshold was originally set quite
243 high for both IWW and EWW samples, but it was later found that IWW gave fewer peaks,
244 possibly due to the complexity of these samples and stronger matrix effects, mostly
245 leading to ionization suppression. The threshold was thus reduced to one quarter to
246 account for this. The maximum number of components was raised to 10,000 to ensure
247 that no compounds would be missed, leading to more than 5,000 components being
248 detected in the IWW and EWW samples (**Table 1**). These were reduced by including

249 compounds with only a m/z 150-500 and $[M+H]^+ = 1$. The in-house database used by the
250 previous two modules within the “Accurate Mass Identification Parameters” of Sieve was
251 then used to gain a list of potential positives.

252 The ChemSpider database (with 10ppm mass accuracy threshold) was also used within
253 Sieve and was begun after the initial component optimization was completed (**Figure 2**).
254 The threshold was made quite high, for optimal “hidden target” analysis, where the
255 detected peaks should correspond to compounds which are commonly and/or highly used.
256 The peak lists of both IWW and EWW results with all data pertaining to mass error, m/z
257 and intensity were exported into Microsoft Excel. From these lists, several thresholds
258 were set and outlying peaks removed: only compounds between m/z 150-500; only
259 $[M+H]^+$; mass error under 2.0ppm; all “hits” just representing a chemical structure, rather
260 than a generic/known name. This final step is rather pragmatic but makes for a more
261 optimal non-target screening, where the remaining compounds should be the more
262 common and/or highly used, as emphasized by having a high intensity threshold.
263 However, this could lead to some less intense peaks being missed and not noted as a
264 possible emerging contaminant in the environment.

265

266 *3.2 Identification with in-house database*

267 Both programs incorporated an in-house database of more than 200 pharmaceuticals,
268 illicit drugs and metabolites (**Table S3**) to get a list of potential positives. All samples
269 were first processed with MsXerator (modules MPeaks and MSCompare) and Sieve 2.1
270 using the parameters outlined in Section 2.5. **Table 2** shows the compounds detected by
271 each.

272 There was very little difference between the compounds found with Sieve and MS
273 Compare, while MPeaks detected somewhat fewer compounds. This could be due to their
274 apparent uses: MPeaks is for pick-peaking, MS Compare for comparing samples, while
275 Sieve does both, resulting in the latter two have more similar results. The fact that all
276 compounds detected by MPeaks were also found with MS Compare leads to the
277 preferential use of the latter module for screening. However, by optimizing the peak-
278 peaking parameters of MPeaks, specifically the sensitivity value, this module could also
279 be of future use in suspect and/or non-target screening.

280 Two methods were used to remove potential false positive peaks: isotopic pattern (for
281 chlorine- and sulfur-containing species) and retention time prediction. Only three of the
282 above compounds (losartan, sulfamethoxazole and temazepam) had a chlorine or sulfur
283 atom, giving rise to a characteristic isotopic pattern. Extracted ion chromatograms were
284 extracted from the initial full-scan data of the Orbitrap and investigated manually. Both
285 sulfamethoxazole and losartan showed the characteristic isotopic pattern, while
286 temazepam did not. Temazepam was thus considered as a false positive and removed
287 from further investigation.

288 A retention time predictor was made, based on artificial neural networks, as in (Miller et
289 al., 2013; Munro et al., 2015) and in our previous work (Bade et al., 2015a). A retention
290 time window of $\pm 11\%$ of total run time was used to find compounds to focus on, based
291 on the window used in our previous work. Of the 25 potential positives investigated, four
292 were removed using this method (benzocaine, ibuprofen, lincomycin and salbutamol)
293 with predicted retention times between 11.5-16.8 minutes (24-36% of the total run time)
294 away from the experimental times. While only four compounds were removed using this
295 technique, it does simplify the identification process, and provides greater confidence in
296 the compounds remaining.

297 *3.3 Identification with ChemSpider*

298 To make a more comprehensive analysis of the samples, an investigation was made using
299 the ChemSpider database search feature of Sieve (**Figure 2**). The introduction of
300 ChemSpider, while removing many components, had the added complication of isobaric
301 and isomeric compounds, with most distinct m/z values having more than one compound
302 associated, as seen in step 5 of **Table 1**. To further refine this list, the mass error was
303 limited to 2ppm (step 6) and all compounds having a formula-only entry were deleted,
304 leaving just compounds with generic names (step 7), leaving up to 100 components in the
305 samples. The literature was then searched to determine whether or not their detection in
306 wastewater could be expected, leading to approximately 30 components and up to 34
307 isomeric/isobaric compounds in the samples. The literature search was made using the
308 Scopus database, and search terms were the generic name of interest, “HRMS”, “LC”,
309 “environment” and “water”. If there were no suitable papers concerning the generic name
310 of interest, the compound was removed from further investigation. To determine which
311 of the isomeric/isobaric compounds the compound within the sample was, the molecular

312 formula was manually searched on ChemSpider, with the compound having the highest
313 number of references deemed to be the compound of interest. Step 7 and the literature
314 search, while pragmatic, were employed to ensure that the compounds detected were
315 those of high consumption/prescription and could therefore be more easily identified in
316 the later *in silico* fragmentation comparison. Finally, eighteen (including three isomers)
317 and eight compounds were finally deemed as potential positives using this non-target
318 approach for IWW and EWW samples respectively (**Table S4**).

319 It is worth noting that by using this approach, most of the same compounds were found
320 as with the in-house database (**Table 2** and **S4**). With such great similarities between the
321 set of potential positive compounds, only 2-hydroxy carbamazepine, desvenlafaxine,
322 adenosine, albendazole, phenazone and the three isomers theophylline, paraxanthine and
323 theobromine required further investigation. Albendazole was the only compound that
324 required an investigation of isotopic pattern as it contains one sulfur atom, which was
325 inconsistent with the mass spectrum, leading to its removal as a false positive. The
326 remaining seven compounds were subjected to retention time prediction based on the time
327 given by Sieve 2.1, and all were found within the set $\pm 11\%$ of total run time retention
328 time window.

329 *3.4 Tentative identification*

330 The potential positives found using both hidden target screening approaches were
331 combined, less those removed in previous steps, leaving 28 compounds to investigate
332 (**Table S1**). These compounds were added to a target list and several IWW and EWW
333 samples were reinjected to see if fragment ions from these compounds could give further
334 confidence to their identification. Metfusion and MassBank were used to help provide
335 further confidence to the fragment ions. As has been mentioned in previous suspect and
336 non-target studies (Agüera et al., 2013; Herrera-Lopez et al., 2014; Zedda and Zwiener,
337 2012), the use and improvement of mass spectral databases, such as MassBank, is
338 extremely important in the tentative identification of compounds for which standards are
339 unavailable. In the end, 22 compounds were able to be tentatively identified (**Table 3**)
340 with at least one fragment ion, while the other six were removed as false positives due to
341 having incorrect fragment ions.

342 One interesting finding was the detection of three isomers (paraxanthine, theobromine
343 and theophylline). These three isomers are all metabolites of caffeine, accounting for

344 80%, 11% and 4% of total metabolism, respectively (Miners and Birkett, 1996).
345 Conventionally, isomeric compounds, separated chromatographically, would be
346 distinguished by retention time. However, as no standards were available, the best way to
347 order the peaks was with retention time prediction. The approach outlined in Sieve 2.1
348 combines the power of the two modules from MsXelerator as described above. After an
349 initial peak-picking process using the settings described in Section 2.5, it compares MS
350 spectra of the procedural blank samples and the studied wastewater samples. Only
351 compounds with an m/z between 150 and 500, and only protonated molecules ($[M+H]^+$)
352 were considered. The list of potential positives were then search by either the in-house
353 database or ChemSpider, which were incorporated in the software. The results of the
354 ChemSpider search were exported into Microsoft Excel, and positive “hits” were
355 considered based on their mass error (< 2 ppm), and if the compound commercially
356 existed i.e. hits which only represented chemical formula were excluded. predicted
357 retention times of 8.48 min, 9.34 min and 9.41 min for theobromine, paraxanthine and
358 theophylline, respectively. While these times are 1-2 minutes from the experimental
359 retention time, they do provide an idea for the order of the isomers. To give more
360 confidence to this information, the fragment ions were checked. As seen in **Figure 3**, the
361 peak at 6.27 min had fragment ions of m/z 163 and 138 while the peaks at 7.98 and 8.37
362 both had one major peak of m/z 124. These fragment ions were checked and compared
363 with MassBank and the literature (Bianco et al., 2009; Gómez et al., 2010; Horai et al.,
364 2010).. Theobromine was found to have fragment ions of m/z 163 and 138, while both
365 paraxanthine and theophylline were found to have a main fragment ion of m/z 124. The
366 losses leading to each fragment ion is defined in **Figure 3**. To differentiate the latter two,
367 the initial retention time predictions led to paraxanthine being the larger peak at 7.98 min
368 and theophylline the small peak at 8.37 min.

369 While 22 compounds were tentatively identified using the workflow outlined throughout
370 this paper, it must be noted that even incorporating the false positive removal strategies
371 of retention time prediction and isotopic pattern as well as fragment ions, the final
372 confirmation of the identity of compounds requires the use of reference standards.
373 Nevertheless, the addition of advanced deconvolution tools (MsXelerator and Sieve) to
374 the HRMS data of the Orbitrap has been shown to be of great value, and the results show
375 how far one can go without the need to purchase reference standards. The information
376 obtained with this strategy circumvents the cost and problems associated with the storage

377 and expiry dates of standards in the laboratories, as the purchase can be directed only
378 towards those compounds that have been previously tentatively identified in the samples.

379

380 **Conclusion**

381 This work has shown that, following initial unbiased, non-target oriented deconvolution
382 using two tools (MsXelerator and Sieve), allowed relevant peaks of interest to be attained.
383 The complementary use of an in-house database or ChemSpider facilitated detection and
384 enabled the identification of more compounds than using just one of these databases.

385 The combination of deconvolution tools and high resolution mass spectrometry, without
386 the use of any reference standards, has enabled 22 compounds to be tentatively identified
387 in environmental water samples. The majority of compounds that were identified in
388 wastewater samples were pharmaceuticals, including the metabolites 4-formylamino
389 antipyrine, 4-acetylamino antipyrine, theobromine, theophylline, paraxanthine and
390 carboxylosartan.

391 It is worth noting that the two hidden target approaches primarily found the same
392 compounds, with some exceptions. Furthermore, when applying small databases it is
393 often easier to analyse the raw data directly. Whereas, when a much larger database is
394 incorporated, these software tools will facilitate searching as well as reducing processing
395 time. With further improvements to these computational programs non-target analysis
396 will become more enticing and easier for laboratories to use in everyday screening
397 methods.

398

399 **Acknowledgements**

400 Richard Bade and Ana Causanilles acknowledge the European Union for their Early Stage
401 Researcher (ESR) contracts as part of the EU-International Training Network SEWPROF
402 (Marie Curie- PEOPLE Grant #317205)

403 Part of this work was supported by the COST Action ES1307 “SCORE - Sewage
404 biomarker analysis for community health assessment”.

405 The financial support of Generalitat Valenciana (Prometeo II 2014/023) and of the
406 Spanish Ministry of Economy and Competitiveness (Project ref CTQ2015-65603) is also
407 acknowledged by the authors of University Jaume I.

408 **References**

- 409 Aceña, J., Stampachiachiere, S., Pérez, S., Barceló, D., 2015. Advances in liquid
410 chromatography–high-resolution mass spectrometry for quantitative and
411 qualitative environmental analysis. *Anal. Bioanal. Chem.* 6289–6299.
412 doi:10.1007/s00216-015-8852-6
- 413 Agüera, A., Martínez Bueno, M.J., Fernández-Alba, A.R., 2013. New trends in the
414 analytical determination of emerging contaminants and their transformation
415 products in environmental waters. *Environ. Sci. Pollut. Res. Int.* 20, 3496–515.
416 doi:10.1007/s11356-013-1586-0
- 417 Bade, R., Bijlsma, L., Miller, T.H., Barron, L.P., Sancho, J.V., Hernández, F., 2015a.
418 Suspect screening of large numbers of emerging contaminants in environmental
419 waters using artificial neural networks for chromatographic retention time
420 prediction and high resolution mass spectrometry data analysis. *Sci. Total Environ.*
421 538, 934–941. doi:10.1016/j.scitotenv.2015.08.078
- 422 Bade, R., Rousis, N.I., Bijlsma, L., Gracia-Lor, E., Castiglioni, S., Sancho, J. V.,
423 Hernandez, F., 2015b. Screening of pharmaceuticals and illicit drugs in wastewater
424 and surface waters of Spain and Italy by high resolution mass spectrometry using
425 UHPLC-QTOF MS and LC-LTQ-Orbitrap MS. *Anal. Bioanal. Chem.* 407, 8979–
426 8988. doi:10.1007/s00216-015-9063-x
- 427 Bianco, G., Abate, S., Labella, C., Cataldi, T.R.I., 2009. Identification and
428 fragmentation pathways of caffeine metabolites in urine samples via liquid
429 chromatography with positive electrospray ionization coupled to a hybrid
430 quadrupole linear ion trap (LTQ) and Fourier transform ion cyclotron resonance
431 mass spec. *Rapid Commun. Mass Spectrom.* 23, 1065–1074.
432 doi:10.1002/rcm.3969
- 433 Bijlsma, L., Emke, E., Hernandez, F., de Voogt, P., 2012. Investigation of drugs of
434 abuse and relevant metabolites in Dutch sewage water by liquid chromatography
435 coupled to high resolution mass spectrometry. *Chemosphere* 89, 1399–1406.
436 doi:10.1016/j.chemosphere.2012.05.110
- 437 Bijlsma, L., Emke, E., Hernández, F., de Voogt, P., 2013. Performance of the linear ion
438 trap Orbitrap mass analyzer for qualitative and quantitative analysis of drugs of
439 abuse and relevant metabolites in sewage water. *Anal. Chim. Acta* 768, 102–110.
440 doi:10.1016/j.aca.2013.01.010
- 441 Bletsou, A.A., Jeon, J., Hollender, J., Archontaki, E., Thomaidis, N.S., 2015. Targeted
442 and non-targeted liquid chromatography-mass spectrometric workflows for
443 identification of transformation products of emerging pollutants in the aquatic
444 environment. *Trends Anal. Chem.* 66, 32–44.
- 445 de Voogt, P., Emke, E., Helmus, R., Panteliadis, P., van Leerdam, J.A., 2011.
446 Determination of illicit drugs in the water cycle by LC-Orbitrap MS, in:
447 Castiglioni, S., Zuccato, E., Fanelli, R. (Eds.), *Illicit Drugs in the Environment:
448 Occurrence, Analysis, and Fate Using Mass Spectrometry*. John Wiley & Sons,
449 Ltd., pp. 87–114.
- 450 Fatta-Kassinos, D., Meric, S., Nikolaou, A., 2011. Pharmaceutical residues in
451 environmental waters and wastewater: Current state of knowledge and future
452 research. *Anal. Bioanal. Chem.* 399, 251–275. doi:10.1007/s00216-010-4300-9

- 453 Gago-Ferrero, P., Schymanski, E.L., Bletsou, A.A., Aalizadeh, R., Hollender, J.,
454 Thomaidis, N.S., 2015. Extended Suspect and Non-Target Strategies to
455 Characterize Emerging Polar Organic Contaminants in Raw Wastewater with LC-
456 HRMS/MS. *Environ. Sci. Technol.* 49, 12333–12341. doi:10.1021/acs.est.5b03454
- 457 Gerlich, M., Neumann, S., 2013. MetFusion: integration of compound identification
458 strategies. *J. Mass Spectrom.* 48, 291–8. doi:10.1002/jms.3123
- 459 Gómez, M.J., Gómez-Ramos, M.M., Malato, O., Mezcua, M., Fernández-Alba, A.R.,
460 2010. Rapid automated screening, identification and quantification of organic
461 micro-contaminants and their main transformation products in wastewater and
462 river waters using liquid chromatography-quadrupole-time-of-flight mass
463 spectrometry with an accurate-mass. *J. Chromatogr. A* 1217, 7038–54.
464 doi:10.1016/j.chroma.2010.08.070
- 465 Helbling, D.E., Hollender, J., Kohler, H.-P.E., Singer, H., Fenner, K., 2010. High-
466 throughput identification of microbial transformation products of organic
467 micropollutants. *Environ. Sci. Technol.* 44, 6621–7. doi:10.1021/es100970m
- 468 Hernández, F., Ibáñez, M., Bade, R., Bijlsma, L., Sancho, J.V., 2014. Investigation of
469 pharmaceuticals and illicit drugs in waters by liquid chromatography-high-
470 resolution mass spectrometry. *TrAC Trends Anal. Chem.* 63, 140–157.
471 doi:10.1016/j.trac.2014.08.003
- 472 Hernández, F., Ibáñez, M., Botero-Coy, A.-M., Bade, R., Bustos-López, M.C., Rincón,
473 J., Moncayo, A., Bijlsma, L., 2015a. LC-QTOF MS screening of more than 1,000
474 licit and illicit drugs and their metabolites in wastewater and surface waters from
475 the area of Bogotá, Colombia. *Anal. Bioanal. Chem.* 407, 6405–6416.
476 doi:10.1007/s00216-015-8796-x
- 477 Hernández, F., Ibáñez, M., Portolés, T., Cervera, M.I., Sancho, J. V, López, F.J., 2015b.
478 Advancing towards universal screening for organic pollutants in waters. *J. Hazard.*
479 *Mater.* 282, 86–95. doi:10.1016/j.jhazmat.2014.08.006
- 480 Hernández, F., Pozo, Ó.J., Sancho, J. V., López, F.J., Marín, J.M., Ibáñez, M., 2005.
481 Strategies for quantification and confirmation of multi-class polar pesticides and
482 transformation products in water by LC–MS2 using triple quadrupole and hybrid
483 quadrupole time-of-flight analyzers. *TrAC Trends Anal. Chem.* 24, 596–612.
484 doi:10.1016/j.trac.2005.04.007
- 485 Hernández, F., Sancho, J. V, Ibáñez, M., Abad, E., Portolés, T., Mattioli, L., 2012.
486 Current use of high-resolution mass spectrometry in the environmental sciences.
487 *Anal. Bioanal. Chem.* 403, 1251–64. doi:10.1007/s00216-012-5844-7
- 488 Herrera-Lopez, S., Hernando, M.D., García-Calvo, E., Fernández-Alba, a. R.,
489 Ulaszewska, M.M., 2014. Simultaneous screening of targeted and non-targeted
490 contaminants using an LC-QTOF-MS system and automated MS/MS library
491 searching. *J. Mass Spectrom.* 49, 878–893. doi:10.1002/jms.3428
- 492 Hogenboom, A.C., van Leerdam, J.A., de Voogt, P., 2009. Accurate mass screening and
493 identification of emerging contaminants in environmental samples by liquid
494 chromatography–hybrid linear ion trap Orbitrap mass spectrometry. *J. Chromatogr.*
495 *A* 1216, 510–519. doi:10.1016/j.chroma.2008.08.053
- 496 Horai, H., Arita, M., Kanaya, S., Nihei, Y., Ikeda, T., Suwa, K., Ojima, Y., Tanaka, K.,
497 Tanaka, S., Aoshima, K., Oda, Y., Kakazu, Y., Kusano, M., Tohge, T., Matsuda,

- 498 F., Sawada, Y., Hirai, M.Y., Nakanishi, H., Ikeda, K., Akimoto, N., Maoka, T.,
499 Takahashi, H., Ara, T., Sakurai, N., Suzuki, H., Shibata, D., Neumann, S., Iida, T.,
500 Tanaka, K., Funatsu, K., Matsuura, F., Soga, T., Taguchi, R., Saito, K., Nishioka,
501 T., 2010. MassBank: a public repository for sharing mass spectral data for life
502 sciences. *J. Mass Spectrom.* 45, 703–14. doi:10.1002/jms.1777
- 503 Hug, C., Ulrich, N., Schulze, T., Brack, W., Krauss, M., 2014. Identification of novel
504 micropollutants in wastewater by a combination of suspect and nontarget
505 screening. *Environ. Pollut.* 184, 25–32. doi:10.1016/j.envpol.2013.07.048
- 506 Kern, S., Fenner, K., Singer, H.P., Schwarzenbach, R.P., Hollender, J., 2009.
507 Identification of Transformation Products of Organic Contaminants in Natural
508 Waters by Computer-Aided Prediction and High-Resolution Mass Spectrometry.
509 *Environ. Sci. Technol.* 43, 7039–7046. doi:10.1021/es901979h
- 510 Kind, T., Fiehn, O., 2007. Seven Golden Rules for heuristic filtering of molecular
511 formulas obtained by accurate mass spectrometry. *BMC Bioinformatics* 8, 105.
512 doi:10.1186/1471-2105-8-105
- 513 Krauss, M., Singer, H., Hollender, J., 2010. LC-high resolution MS in environmental
514 analysis: from target screening to the identification of unknowns. *Anal. Bioanal.*
515 *Chem.* 397, 943–951. doi:10.1007/s00216-010-3608-9
- 516 Leendert, V., Van Langenhove, H., Demeestere, K., 2015. Trends in liquid
517 chromatography coupled to high-resolution mass spectrometry for multi-residue
518 analysis of organic micropollutants in aquatic environments. *TrAC Trends Anal.*
519 *Chem.* 67, 192–208. doi:10.1016/j.trac.2015.01.010
- 520 Letzel, T., Bayer, A., Schulz, W., Heermann, A., Lucke, T., Greco, G., Grosse, S.,
521 Schüssler, W., Sengl, M., Letzel, M., 2015. LC – MS screening techniques for
522 wastewater analysis and analytical data handling strategies : Sartans and their
523 transformation products as an example. *Chemosphere* 137, 198–206.
524 doi:10.1016/j.chemosphere.2015.06.083
- 525 Little, J.L., Williams, A.J., Pshenichnov, A., Tkachenko, V., 2012. Identification of
526 “Known Unknowns” Utilizing Accurate Mass Data and ChemSpider. *J. Am. Soc.*
527 *Mass Spectrom.* 23, 179–185. doi:10.1007/s13361-011-0265-y
- 528 Luo, Y., Guo, W., Ngo, H.H., Nghiem, L.D., Hai, F.I., Zhang, J., Liang, S., Wang, X.C.,
529 2014. A review on the occurrence of micropollutants in the aquatic environment
530 and their fate and removal during wastewater treatment. *Sci. Total Environ.* 473-
531 474, 619–41. doi:10.1016/j.scitotenv.2013.12.065
- 532 Makarov, A., Scigelova, M., 2010. Coupling liquid chromatography to Orbitrap mass
533 spectrometry. *J. Chromatogr. A* 1217, 3938–45. doi:10.1016/j.chroma.2010.02.022
- 534 Miller, T.H., Musenga, A., Cowan, D.A., Barron, L.P., 2013. Prediction of
535 chromatographic retention time in high-resolution anti-doping screening data using
536 artificial neural networks. *Anal. Chem.* 85, 10330–7. doi:10.1021/ac4024878
- 537 Miners, J.O., Birkett, D.J., 1996. The use of caffeine as a metabolic probe for human
538 drug metabolizing enzymes. *Gen. Pharmacol.* 27, 245–249. doi:10.1016/0306-
539 3623(95)02014-4
- 540 Munro, K., Miller, T.H., Martins, C.P.B., Edge, A.M., Cowan, D. a., Barron, L.P., 2015.
541 Artificial neural network modelling of pharmaceutical residue retention times in

542 wastewater extracts using gradient liquid chromatography-high resolution mass
543 spectrometry data. *J. Chromatogr. A* 1396, 33–44.
544 doi:10.1016/j.chroma.2015.03.063

545 Schollée, J.E., Schymanski, E.L., Avak, S.E., Loos, M., Hollender, J., 2015. Prioritizing
546 Unknown Transformation Products from Biologically-Treated Wastewater using
547 High-Resolution Mass Spectrometry, Multivariate Statistics, and Metabolic Logic.
548 *Anal. Chem.* acs.analchem.5b02905. doi:10.1021/acs.analchem.5b02905

549 Schymanski, E.L., Jeon, J., Gulde, R., Fenner, K., Ruff, M., Singer, H.P., Hollender, J.,
550 2014a. Identifying Small Molecules via High Resolution Mass Spectrometry:
551 Communicating Confidence. *Environ. Sci. Technol.* 48, 2097–2098.
552 doi:10.1021/es5002105

553 Schymanski, E.L., Singer, H.P., Longrée, P., Loos, M., Ruff, M., Stravs, M.A., Ripollés
554 Vidal, C., Hollender, J., 2014b. Strategies to characterize polar organic
555 contamination in wastewater: exploring the capability of high resolution mass
556 spectrometry. *Environ. Sci. Technol.* 48, 1811–8. doi:10.1021/es4044374

557 van der Aa, M., Bijlsma, L., Emke, E., Dijkman, E., van Nuijs, A.L.N., van de Ven, B.,
558 Hernández, F., Versteegh, A., de Voogt, P., 2013. Risk assessment for drugs of
559 abuse in the Dutch watercycle. *Water Res.* 47, 1848–57.
560 doi:10.1016/j.watres.2013.01.013

561 Zedda, M., Zwiener, C., 2012. Is nontarget screening of emerging contaminants by LC-
562 HRMS successful? A plea for compound libraries and computer tools. - *Anal.*
563 *Bioanal. Chem.* 2493–2502. doi:10.1007/s00216-012-5893-y

564

565

566 **Table 1:** *Number of components after each step of the Sieve hidden target identification approaches..*

567

568

Step	IWW components	EWW components
1.	6690	5091
2.	5158	3528
3.	2014	2175
4.	18	16
NON-TARGET SCREENING		
Number of distinct <i>m/z</i> (total number of compounds)		
5.	239 (437)	441 (677)
6.	166 (362)	308 (543)
7.	100 (150)	64 (108)
Final	18	8

569

570

571 **Table 2:** All compounds detected (in at least one sample) by each program following suspect screening.

Compound	IWW			EWW		
	Sieve	MSCompare	Mpeaks	Sieve	MSCompare	Mpeaks
4-acetylamino antipyrine	Not detected	Detected	Detected	Detected	Detected	Detected
4-formylamino-antipyrine	Not detected	Not detected	Not detected	Detected	Detected	Not detected
Acetaminophen	Detected	Detected	Detected	Detected	Detected	Detected
Benzocaine	Detected	Detected	Detected	Not detected	Detected	Not detected
Benzoylecgonine	Detected	Detected	Detected	Detected	Detected	Detected
Caffeine	Detected	Detected	Detected	Detected	Detected	Detected
Carbamazepine	Not detected	Not detected	Not detected	Detected	Detected	Detected
Cocaine	Detected	Detected	Detected	Not detected	Not detected	Not detected
Cotinine	Detected	Detected	Not detected	Detected	Detected	Not detected
Gemfibrozil	Detected	Not detected	Not detected	Not detected	Not detected	Not detected
Ibuprofen	Not detected	Not detected	Not detected	Detected	Detected	Detected
Irbesartan	Detected	Detected	Not detected	Detected	Detected	Detected
Ketoprofen	Detected	Detected	Not detected	Detected	Detected	Detected
Lidocaine	Detected	Detected	Detected	Not detected	Not detected	Not detected
Lincomycin	Detected	Detected	Detected	Detected	Detected	Detected
Losartan	Detected	Detected	Detected	Detected	Detected	Detected
Metoprolol	Detected	Detected	Detected	Not detected	Not detected	Not detected
Naproxen	Detected	Not detected	Not detected	Detected	Detected	Detected
Phenacetin	Not detected	Not detected	Not detected	Detected	Detected	Detected
Phenytoin	Detected	Not detected	Not detected	Not detected	Detected	Not detected
Salbutamol	Not detected	Detected	Detected	Not detected	Detected	Detected
Sulfamethoxazole	Detected	Detected	Detected	Detected	Detected	Detected
Temazepam	Not detected	Detected	Not detected	Not detected	Not detected	Not detected
Trimethoprim	Detected	Detected	Detected	Not detected	Detected	Not detected
Valsartan	Detected	Detected	Not detected	Detected	Detected	Detected

572

	Detected by the program
	Not detected by the program

573

574

575 **Table 3:** All compounds tentatively identified (level 2a), together with retention time and fragment ions

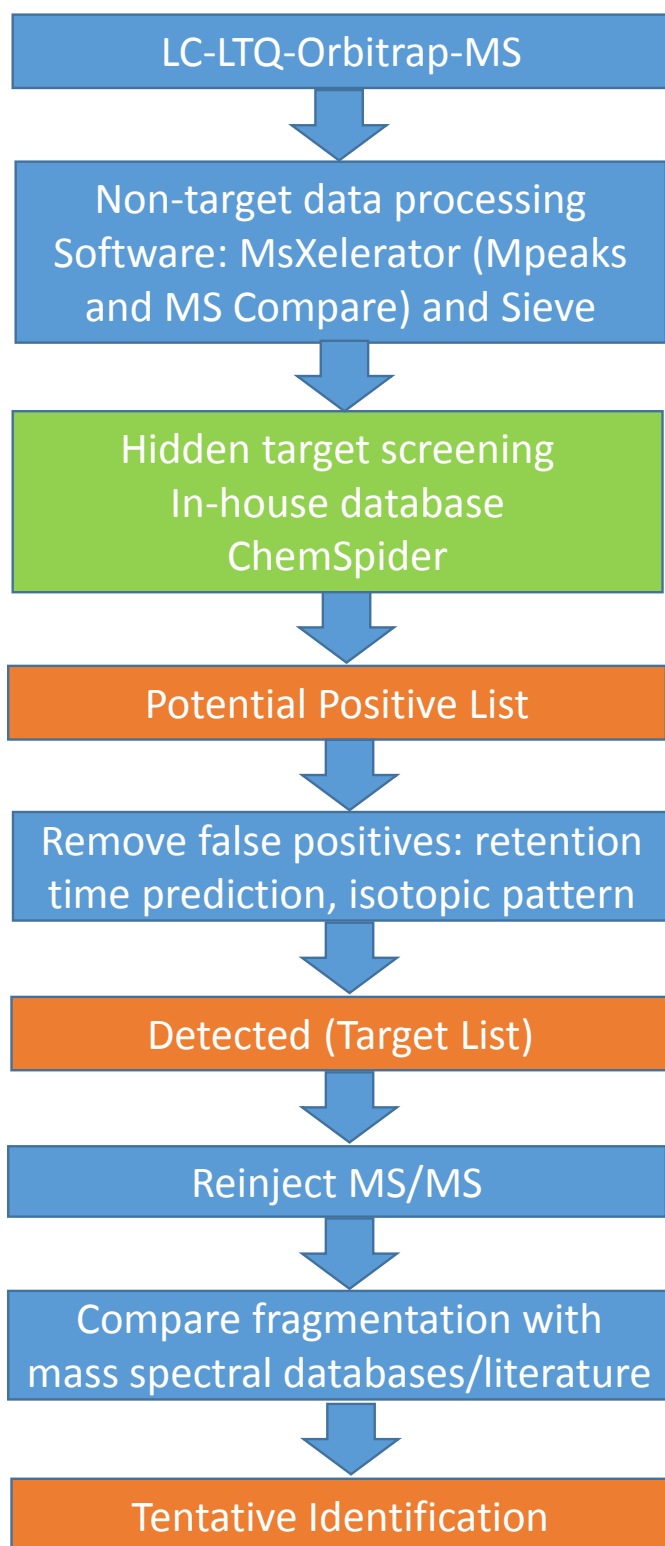
Compound	m/z	RT	Fragment ions			IWW	EWW
4-acetylaminoantipyrine¹	246.1234	9.43	228.1	204.1		X	X
4-formylaminoantipyrine²	232.1086	9.29	214.1	204.1		X	X
Acetaminophen³	152.0706	6.09	110.1	134.0		X	X
Adenosine²	268.1035	3.58	136.0				X
Benzoylecgonine²	290.1385	12.22	168.1			X	X
Caffeine¹	195.0876	10.4	138.1			X	X
Carbamazepine³	237.1022	22.56	194.1	152.9		X	X
Carboxylosartan¹	437.1480	26.94	207.1	235.1	365.3	X	X
Cocaine²	304.1543	13.84	182.2			X	
Ketoprofen³	255.1014	16.78	237.1	209.1		X	X
Lidocaine¹	235.1807	10.12	86.1			X	
Losartan¹	423.1695	25.58	405.0	207.2	377.2	X	X
Metoprolol²	268.1908	13.71	218.1	191.1	159.1	X	
Naproxen³	231.1016	27.24	185.1			X	X
Paraxanthine²	181.0721	7.98	124.1			X	X
Phenacetin²	180.1030	17.22	138.1	110.0		X	X
Phenazone²	189.1022	12.15	161.2	146.1	131.1	X	X
Sulfamethoxazole²	254.0594	12.41	235.8	188.1	156.1	X	X
Theobromine^{2,4}	181.0721	6.28	163.1	137.1	138.1	X	X
Theophylline^{2,4}	181.0721	8.37	124.1			X	X
Trimethoprim²	291.1454	9.91	230.2	123.2	261.1	X	X
Valsartan²	436.2341	28.96	335.1	265.2	155.1	X	X

576 RT = retention time (minutes)

577 The parent ions were recorded at accurate mass (full-scan mode) while the fragment
578 ions were recorded as part of a product ion scan in the ion trap part with nominal mass
579 measurement.580 ¹ Information on fragment ions from Hernández *et al* (Hernández *et al.*, 2015a)581 ² Information on fragment ions from MassBank (Horai *et al.*, 2010)582 ³ Information on fragment ions from Bade *et al* (Bade *et al.*, 2015b)583 ⁴ Information on fragment ions from Gómez *et al* (Gómez *et al.*, 2010)

584

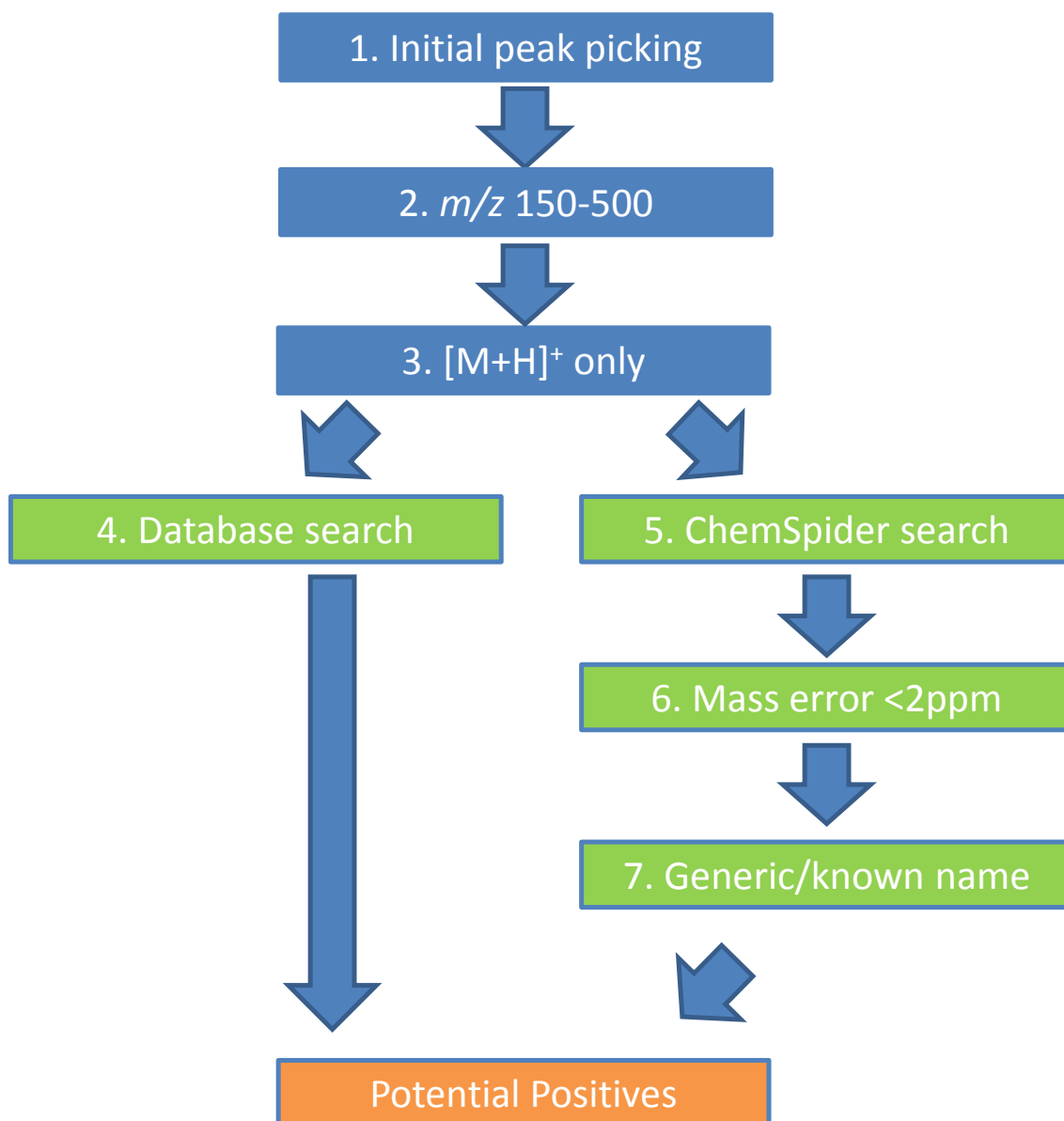
585



586

587 **Figure 1:** Workflow for screenings using the deconvolution tools MsXelerator and Sieve. All orange levels
588 represent specific identification confidence levels

589



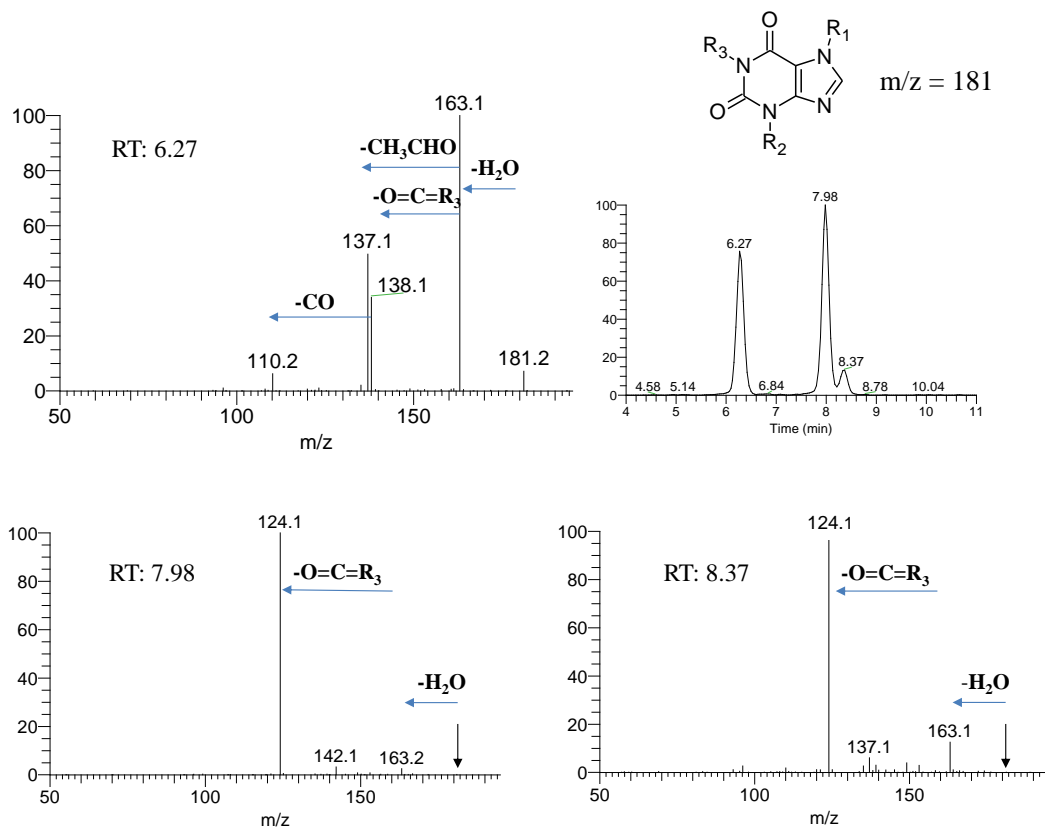
590

591 **Figure 2:** Sieve workflow for the two hidden target identification approaches.

592

593

594



595

596

597

598

599

Figure 3: Tentative identification of theobromine (top left), paraxanthine (bottom left) and theophylline (bottom right), with chromatographic peaks (top right). The generic structure has been shown in the top right corner, where R_1 , R_2 and R_3 differ for the metabolites as follows: theobromine: $R_1 = CH_3$, $R_2 = CH_3$ and $R_3 = H$; theophylline: $R_1 = H$, $R_2 = CH_3$ and $R_3 = CH_3$; paraxanthine: $R_1 = CH_3$, $R_2 = H$ and $R_3 = CH_3$

600