



Stockholm  
University

Department of Materials- and Environmental Chemistry  
Unit of Analytical Chemistry

**Widening the Horizon of Non-Targeted  
Environmental LC/ESI/HRMS Analysis:  
Preparation for the NORMAN  
Interlaboratory Comparison on Semi-  
Quantification**

15 Credits  
Bachelor's Thesis

Louise Malm

Supervisor: Anneli Kruve, Associate Professor  
Stockholm, July 2020

# Table of Contents

<b>Table of Contents .....</b>	<b>2</b>
<b>Abbreviations .....</b>	<b>3</b>
<b>Abstract.....</b>	<b>4</b>
<b>1. Introduction.....</b>	<b>5</b>
<b>1.1. Target, Suspect and Non-Targeted Screening .....</b>	<b>5</b>
<b>1.2. NORMAN Network.....</b>	<b>7</b>
<b>1.3. Interlaboratory Comparisons.....</b>	<b>8</b>
<b>1.4. Response factor .....</b>	<b>8</b>
<b>1.5. Strategies for Semi-Quantification .....</b>	<b>9</b>
<b>1.5.1. Structurally Similar Compounds .....</b>	<b>9</b>
<b>1.5.2. Parent Compound – TPs .....</b>	<b>9</b>
<b>1.5.3. Close Eluting Compounds .....</b>	<b>9</b>
<b>1.5.4. Predicting the Ionisation Efficiency.....</b>	<b>10</b>
<b>2. Materials and Methods.....</b>	<b>11</b>
<b>2.1. Chemicals.....</b>	<b>11</b>
<b>2.2. Samples .....</b>	<b>11</b>
<b>2.3. Instrumental.....</b>	<b>12</b>
<b>2.4. Data Treatment.....</b>	<b>12</b>
<b>3. Results.....</b>	<b>13</b>
<b>4. Discussion .....</b>	<b>16</b>
<b>4.1. Compound Selection and Distribution .....</b>	<b>16</b>
<b>4.2. Prediction Errors .....</b>	<b>17</b>
<b>4.2.1. Correlation of the Errors with Chromatographic Properties .....</b>	<b>18</b>
<b>4.2.2. Adduct Formation and Semi-Quantification Error.....</b>	<b>19</b>
<b>4.3. Structural Similarity .....</b>	<b>19</b>
<b>4.3.1. Parent Compound Approach .....</b>	<b>20</b>
<b>4.3.2. Similar Compound Approach.....</b>	<b>20</b>
<b>4.3.3. Closest Eluting Standard Approach .....</b>	<b>21</b>
<b>4.3.4. Ionisation Efficiency Approach.....</b>	<b>24</b>
<b>4.4. Stability .....</b>	<b>25</b>
<b>5. Conclusion .....</b>	<b>26</b>
<b>6. Future Challenges .....</b>	<b>27</b>
<b>7. Acknowledgements .....</b>	<b>28</b>
<b>8. References .....</b>	<b>29</b>
<b>9. Supplementary Information.....</b>	<b>31</b>

## Abbreviations

DMSO	Dimethyl sulfoxide
ESI	Electrospray ionisation
HRMS	High resolution mass spectrometer
<i>IE</i>	Ionisation efficiency
ILIS	Isotope labelled internal standards
LC	Liquid chromatography
MeCN	Acetonitrile
MeOH	Methanol
RF	Response factor
RT	Retention time
SI	Supplementary information
SPE	Solid phase extraction
TPs	Transformation products

## **Abstract**

The number of chemicals in daily use increases and surface water worldwide is polluted by an ever-widening spectrum of compounds, from pharmaceutical residues and beauty products to pesticides and artificial sweeteners. The quality of the surface water influences the purity of our drinking water and also puts higher demands on wastewater treatment plants. The NORMAN Network<sup>1</sup> has compiled a list of more than 65 000 compounds suspected to be found in surface water. Detecting all of these compounds requires various analytical techniques, of which liquid chromatography electrospray ionisation high resolution mass spectrometry (LC/ESI/HRMS) is the preferred technique for water analysis. However, ordinary targeted analysis and quantification with standard substances is neither time nor cost effective for so many analytes. Furthermore, such analysis cannot detect new, emerging pollutants or ones for which there are no standards available. Therefore, non-targeted analysis is being increasingly employed. Together with semi-quantification of the analytes of interest, it is possible to prioritise which analytes are most important to identify.<sup>2,3</sup> Though there are semi-quantification methods available for use, none are yet standardised, and therefore not ready for wide-community use.

The NORMAN network is arranging a large interlaboratory trial for semi-quantitative non-targeted environmental analysis in 2020 to 2021. The intention of this trial is to make semi-quantification results comparable between different instruments and laboratories, but also to validate the proposed semi-quantification methods. The current thesis is a preparation for the collaborative trial, with focus on selecting the compounds to be used in the large-scale comparison. Furthermore, four semi-quantification approaches will be tested and evaluated.

# 1. Introduction

Even though our planet is largely covered in water, access to freshwater in some areas is limited and the scarce quantities available is steadily contaminated. Water pollutants from an ever growing spectrum of contaminants is not only a threat to global health in the sense of increased shortage of drinking water and sanitation,<sup>4</sup> but it also threatens life in the waters as well as the animals that eat for instance contaminated fish.<sup>5</sup> Many contaminants, up to approximately 70 000 according to Schwarzenbach et al. (2006),<sup>6</sup> come from common household products used daily such as sunscreens, medicines, detergents, etc. via the sewage systems and insufficient treatment plants. Agriculture with associated pesticides and fertilisers, together with industry are, however, responsible for the majority of pollutants, releasing hundreds of millions of tons to the environment annually.<sup>6</sup> Many of these compounds escape through wastewater treatment plants and end up in the surface water. Additionally, many chemicals degrade, either naturally in the environment or during water purification,<sup>7</sup> and thereby give rise to a number of transformation products (TPs) with supposedly similar structure as their parent compound. A comprehensive list of suspected compounds found in water by the NORMAN Network<sup>1</sup> contains more than 65 000 entries, giving an indication of the seriousness of the problem. Many of these compounds are TPs, mainly from pesticides, though contaminants from other groups like pharmaceuticals most likely also undergo degradation.<sup>8,9</sup> Furthermore, TPs are sometimes both more toxic and more abundant than their parent compound,<sup>10</sup> thus they should not be ignored.

Some steps to control the amount of pollutants released to the aquatic environment are in place, e.g. the Water Framework Directive in EU and Clean Water Act in the US. Still, the abovementioned organisations only control 48 compounds<sup>11</sup> and 126 compounds,<sup>12</sup> respectively, i.e. merely a tiny fraction of all compounds currently contaminating the surface waters.<sup>9</sup> Moreover, TPs are seldom included in water protective plans, mainly due to unknown structure<sup>8</sup> or toxicity.<sup>10</sup>

One of the more common techniques to analyse water samples for contaminants is liquid chromatography electrospray ionisation high resolution mass spectrometry (LC/ESI/HRMS) due to the polar nature of many contaminants.<sup>8</sup> Furthermore, a higher diversity of analytes that can be detected in LC/MS compared to gas chromatography mass spectrometry (GC/MS), since the compounds do not have to be volatile.<sup>13</sup> To increase the contaminants controlled, including more TPs, it is not feasible to use ordinary targeted analysis, i.e. to only look for a couple of selected pollutants in the sample. Neither can this kind of analysis find new emerging compounds, since it only finds *what* it's looking for. Instead, non-targeted or suspect screening methods are used more and more. However, a crucial problem remains, namely quantification. It is important to know both *which* pollutants are found in water, as well as *how much* of them there are. A solution to this lies in semi-quantification, which unlike "real" quantification where analytical standards are used, estimates the amount of the analytes without standards. This is an ongoing research topic, where new approaches for semi-quantification are just being, or have very recently been, developed. In this study, four of those strategies will be evaluated, which are further described in section 1.5. below.

## 1.1. Target, Suspect and Non-Targeted Screening

Target screening is a method to look for predetermined compounds, and with the help of reference standard compounds it allows quantification of the analytes. This is the default method to screen for controlled contaminants in water; however, this approach is not sufficient to screen for all contaminants that can be found in surface water.

In order to find more, and above all new, contaminants in water, non-targeted screening needs to be employed. This procedure is much more time consuming than targeted screening, especially the data processing part.

For identification, the chromatographic peaks first need to be sorted and combined according to which species they belong to.<sup>8</sup> This is done by a software, usually from the same manufacturer as the hardware, e.g. Compound Discoverer (Thermo Fisher Scientific™, USA). However, free software can also be used to combine the chromatographic peaks, e.g. from <https://www.rformassspectrometry.org/>. This step gives a set of possible molecular formulae for each peak, and after additional filtering, e.g. matching isotope pattern, one molecular formula is assigned to each peak. Identification by molecular formula reaches confidence level 4 according to a scale proposed by Schymanski et al. (2014),<sup>14</sup> as seen in figure 1.

To obtain a tentative structure of the contaminant, further processing is needed, usually by analysing the fragmentation in MS<sup>2</sup> spectra. In the rare case that only one tentative structure is obtained from MS<sup>2</sup> experiments, and additional evidence support that structure, confidence level 2 might be achieved. If instead more than one possible structure is procured, the fragmentation spectra needs to be compared to reference spectra to reach level 2 identification. The level is, therefore, called probable structure, indicating that the structure is assigned with rather high confidence, but not as high as the confirmed structure level. To reach the highest confidence level, the probable structure needs to be confirmed with the reference standard.

Example	Identification confidence	Data requirement
	<b>1: Confirmed structure</b> by reference standard	MS, MS <sup>2</sup> , RT, Reference standards
	<b>2: Probable structure</b> a) by library spectrum match b) by diagnostic evidence	MS, MS <sup>2</sup> , Library MS <sup>2</sup> MS, MS <sup>2</sup> , Experimental data
	<b>3: Tentative candidate(s)</b> structure, substituent, class	MS, MS <sup>2</sup> , Experimental data
C <sub>7</sub> H <sub>12</sub> ClN <sub>5</sub>	<b>4: Molecular formula</b> unequivocal	MS isotope/adduct
201.6567	<b>5: Exact mass</b> of interest	MS

**Figure 1.** The confidence levels for identification with their respective data requirements, as proposed by Schymanski et al.<sup>14</sup> As seen, a (tentative) structure is only available from level 3 and higher.

The steps taken from level 3 to full identification usually requires searching databases for a match and is therefore not really “true” non-targeted analysis but rather suspect analysis. The

main difference between suspect- and non-targeted screening are that in suspect screening there is an idea of what to search for, e.g. from suspect lists or databases, whereas non-target analysis is performed with no prior information to what to search for.

Understandably, data treatment in non-targeted and suspect screening is tiresome, and some compounds may be overlooked or wrongfully dismissed in the process. Similarly, steps taken prior to data analysis, e.g. sampling, preparation, separation and detection, can also introduce errors.<sup>8</sup>

As some compounds, especially in water samples, are quite unstable it is generally important to analyse the sample as shortly after sampling as possible. If too long time passes between sampling and analysis, some compounds might degrade and thus cannot be detected. This is especially important for surface water, since this kind of water is fairly “new” as opposed to the much older groundwater. Furthermore, analysing in close proximity to sampling also minimise the risk for contamination.

Most often some kind of sample preparation is required to concentrate the sample, or many compounds would not be detected as their concentration is below instrument sensitivity.<sup>10,15,16</sup> Solid phase extraction (SPE) is a commonly used preconcentration technique, and is also used in this study. The choice of cartridge highly influences which compounds will be detected in the analysis and must therefore be carefully considered. However, since the analytes are unknown in non-targeted analysis, choosing the solid phase can be quite hard. To prevent the removal of analytes of interest, extracts from cartridges of different polarity can be combined. Preferably, a multilayer cartridge can be used, where the different layers have different polarity.<sup>7</sup> A standard cartridge used in non-targeted screening of water samples is OASIS® HBL, since it contains a copolymer with both hydrophilic and hydrophobic moieties.<sup>17</sup>

Another source to false negatives is introduced in the LC/HRMS process, or rather in the ionisation process. Since the mass spectrometer can only detect ions, compound that do not ionise well enough will not be detected, but this does not mean that the compound is not present in the sample. It is also important to consider the choice of column used in LC, to obtain the best chromatographic separation.

## 1.2. NORMAN Network

The NORMAN Association is a “network of reference laboratories, research centres and related organisations for monitoring of emerging environmental substances”, i.e. a network of institutions that in some way deals with emerging substances. It is a non-profit organisation, aiming to collect and maintain information and data of new environmental substances, as well as promote further development of measurement and monitoring tools. The NORMAN Network manages or contributes to many databases, e.g. MassBank Europe (mass spectra of emerging substances) and Suspect List Exchange (various lists of substances for suspect screening and prioritisation), free for anyone to use. Furthermore, they seek to improve and exchange knowledge of analytical methods and approaches, mainly by interlaboratory studies. Collaborating laboratories are from all over Europe and North America.<sup>18</sup>

One of the interlaboratory collaborations arranged by the NORMAN Network this year (2020) is on semi-quantitative non-targeted analysis with LC/ESI/HRMS,<sup>19</sup> with the intention to make semi-quantification strategies comparable between different instruments and laboratories.

### 1.3. Interlaboratory Comparisons

There are different purposes for doing interlaboratory comparisons, which needs to be established before designing such studies. The goal can vary from determining the physical properties of an artefact to determining random or systematic variations in measurement results,<sup>20</sup> or to evaluate either a specific measurement method or a new reference material.<sup>21</sup>

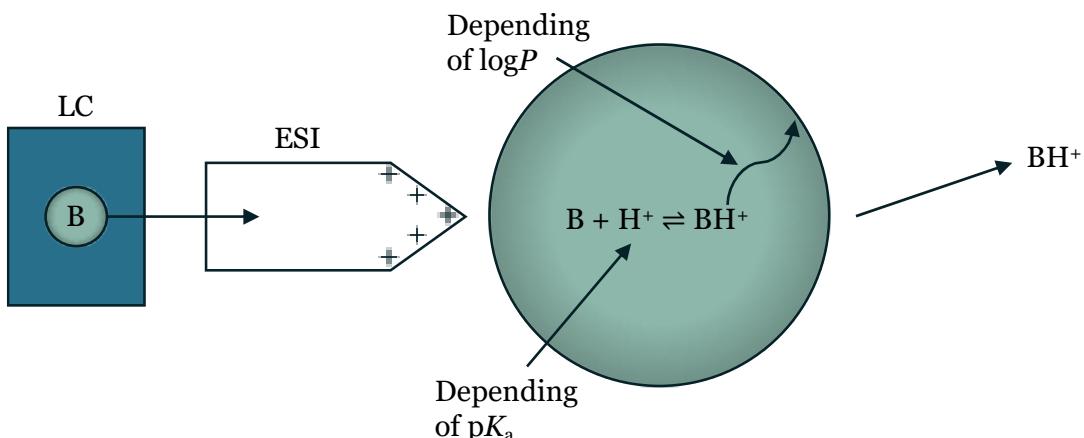
In this particular case, the objective is to validate the semi-quantification approaches currently used and make the results comparable between laboratories. To do this, each participating laboratory will receive an aliquot from the same water sample, spiked with 67 chosen compounds. Of these, 37 will have a, for the laboratory unknown concentration, which will be determined using the quantification strategies.<sup>22</sup> The estimated concentrations will then be compared with the actual concentrations for each of the 37 compounds, to evaluate the semi-quantification approaches. Since the plan for the interlaboratory comparison is to send already spiked water samples to the laboratories, it is important to know whether the substances are stable in water matrix. Therefore, the stability will be tested and evaluated during this pre-study.

### 1.4. Response factor

The response factor (RF) of a compound is the ratio between the detected peak area and the concentration of the compound (see equation 1), with the unit [signal unit/M], where the signal is in arbitrary units.

$$RF = \frac{\text{peak area}}{\text{concentration}} \quad (1)$$

The RF is highly influenced by how well the compound ionise; therefore, since not all compounds ionise to the same extent, the peak area is not proportional to the concentration. This is giving rise to problems regarding quantification and is by extension one of the reasons why semi-quantification methods are needed. Many aspects are influencing the efficiency with which a substance ionises, e.g. the properties of the substance and the eluent, matrix effects and which ionisation source is used.<sup>3,23</sup> Important properties regarding the compound itself is hydrophobicity and basicity/acidity, as seen in figure 2. Regarding the eluent, the amount of organic solvent is significant, as well as the pH and buffer type.<sup>8</sup> Generally, the ionisation efficiency (IE) increases with the amount of organic solvent, since this makes the droplets in the ionisation source dry faster.<sup>3</sup>



**Figure 2.** Two important properties of the compound (B) that is influencing how well it will ionise. The basicity ( $pK_a$ ) is influencing how well the substance will become protonated inside the droplet, and the hydrophobicity ( $\log P$ ) is influencing how easy the protonated ion will get to the surface of the droplet.

## 1.5. Strategies for Semi-Quantification

In the interlaboratory trial, a total of six semi-quantification strategies will be tested and evaluated; however, in this smaller pre-study, only four of those will be examined due to the current technical availability. Namely, two of the proposed approaches are still under development. The used semi-quantification strategies are described in more detail below, with their respective assumptions and limitations displayed.

### 1.5.1. Structurally Similar Compounds

This approach takes advantage of the assumption that compounds with structural similarity have similar RF. This means, that the most structurally similar compound from the standards mix will be used to quantify the compound in question, according to equation 2:

$$c_{\text{suspect compound}} = \frac{\text{peak area}_{\text{suspect compound}}}{\text{RF}_{\text{most similar standard compound}}} \quad (2)$$

The most similar compound is found by comparing the functional groups of the suspect compound with all of the standard compounds, together with some other 2D chemical descriptors like the distance between functional groups.<sup>22</sup> An online tool<sup>24</sup> is available to search for the structurally most similar compound from the University of Athens standards database, down to a similarity of 5%. This tool requires the SMILES string of the suspected compound as an input. In this study, the list of similar compounds found by the tool is compared with the list of compounds in the standards mix to find the most structurally similar standard available. Additionally, the similarity score (%) and a maximum expected error is obtained.

A disadvantage with the online tool is that it calculates the similarity based on the *number* of similar functional groups but does not consider *which* groups are similar or dissimilar. E.g. the most similar compound may lack a very hydrophobic or basic group compared to the suspect compound, and still have a very high similarity score. This disadvantage might make this approach prone to errors.

### 1.5.2. Parent Compound – TPs

This strategy is quite similar to the previous but is instead taking advantage of the fact that TPs originate from a parent compound. The parent compound is then used, similarly as the structurally similar standard, to estimate the concentration of its TPs in the unknown mixture. This approach presumes that TPs are structurally similar to their parent, and thus are assumed to behave similarly. With this assumption, the RF of the parent is directly applicable for the TP, giving the following equation to calculate the concentration of the transformation product:

$$c_{TP} = \frac{\text{peak area}_{TP}}{\text{RF}_{\text{parent compound}}} \quad (3)$$

Dahal et al. (2011) has previously applied this strategy on drugs and their metabolites, with up to a 4-fold error.<sup>25</sup> However, this approach cannot be used as a single semi-quantification method in true non-targeted screening, as it is solely applicable for compounds that are the result of degradation.

### 1.5.3. Close Eluting Compounds

The third semi-quantification approach is based on previous work by Pieke et al. (2017),<sup>2</sup> and uses the response factor of the standard with the closest RT to the suspect compound, as per equation 4:

$$c_{\text{suspect compound}} = \frac{\text{peak area}_{\text{suspect compound}}}{\text{RF}_{\text{closest eluting standard compound}}} \quad (4)$$

Here, the estimated concentration of the suspect compound is not necessarily based on the standard with the most similar structure, but rather on the standard with the most similar chromatographic properties. The assumption is that compounds with similar retention time have resembling *IE* as well, and consequently, the RF of close eluting standard can be applied to the suspect compound. This approach has previously been employed with satisfying results, with a 2-fold average error.<sup>2</sup> Furthermore, this approach can be applied for all compounds detected in a sample, and even for compounds without a tentative structure.<sup>22</sup>

#### *1.5.4. Predicting the Ionisation Efficiency*

Lastly, a somewhat different approach will be tested, where the response factor of the suspect compound is predicted. The prediction is based on a model recently developed by Liigand et al. (2020),<sup>23</sup> which uses 2D PaDEL descriptors of the compound, eluent descriptors and random forest regression. Especially important parameters regarding the eluent are pH, eventual presence of NH<sub>4</sub><sup>+</sup>, as well as viscosity, surface tension and polarity index. The *IE* can be calculated in R (© The R Foundation) based on the previously developed model, or by an online tool provided by Quantem Analytics.<sup>26</sup> The known concentrations of the standards will be used to calibrate the algorithm used for predictions to match the analytical method used.

Similar to some of the abovementioned strategies, this strategy can be used in suspect screening, i.e. when a tentative structure is available. However, a drawback is that currently, the algorithms can only predict the response factor for protonated ions. This means that this approach is not applicable for compounds that form only other major adducts, e.g. sodium or ammonium adducts.

## 2. Materials and Methods

### 2.1. Chemicals

The standard mix consisted of Guanylurea, Amitrole, Histamine, Chlormequat, Methamidophos, Vancomycin, Trichlorfon, Butocarboxim, Dichlorvos, Tylosin, Rifaximin, Spinosyn A, Emamectin B1a, Nigericin, Ivermectin B1a, TCMTB chosen by the University of Athens, and Atrazine, Octocrylene, Clarithromycin, Aspartame, Simvastatin, Sucralose, Igarol, Caffeine, Carbamazepine, Benzotriazole, Metolachlor and Imazalil selected at Stockholm University.

The suspect mix consisted of Atrazine-desethyl, Atrazine-desethyl-2-hydroxy, Atrazine-desisopropyl, Atrazine-desethyl-desisopropyl, Atrazine-desethyl-desisopropyl-2-hydroxy, Atrazine-desisopropyl-2-hydroxy, Atrazine-2-hydroxy, 2-(Methylthio)benzothiazole, Progesterone, Butylamine, Haloperidol, Reserpine, Phenazine, Clotrimazole, Simazine, Efavirenz, Adenosine, Climbazole, Melamine, Metazachlor, Chlorothiazide, Metformin, 2-Methylbenzothiazole, Benzothiazole, Chlorpyrifos, 5-Methyl-1H-benzotriazole, 10,11-Dihydro-10-hydroxycarbamazepine, Sudan I, Ketoconazole, 5-Chlorobenzotriazole, Benzotriazole-5-carboxylic acid, Carbamazepine-10,11-epoxide, Metolachlor-ESA, Metolachlor-OXA, Omethoate, 2-Hydroxybenzothiazole and 2-Aminobenzothiazole.

The isotope labelled internal standard (ILIS) mix consisted of Atrazine-d<sub>5</sub>, Caffeine-<sup>13</sup>C<sub>3</sub> and Imazalil-d<sub>5</sub>.

TCMTB was bought from Honeywell Fluka™, all other chemicals were from Sigma-Aldrich. All chemicals were of analytical standard.

Each chemical was dissolved in either Acetonitrile (MeCN), Methanol (MeOH), MilliQ water, Dimethyl sulfoxide (DMSO), 0.1% formic acid or 0.1M Hydrochloric acid (HCl), or in a mix of two of the mentioned solvents. The exact solvent for each compound can be found in table S1 in supplementary information (SI). All solvents were from Honeywell Riedel-de Haën™, except HCl that were from VWR Chemicals and ultrapure water which was purified with Milli-Q IQ 7000 (Merck, Darmstadt, Germany).

### 2.2. Samples

Solutions of all chemicals were prepared to a concentration of around 1000 ppm, mixed according to the mixes in section 2.1 and were lastly diluted. Final concentration for each compound and mix can be found in tables S1 and S2 in SI.

Water sample from Laduviken was obtained on June 16<sup>th</sup>, filtered twice using Munktell Filter Paper (Ahlstrom Munksjö) and stored at +4°C prior to sample preparation. Tap water samples were obtained from a tap at Stockholm University directly before extraction.

SPE was performed to concentrate the water samples, based on the SPE method used by Rousis et al. (2017),<sup>27</sup> using OASIS® HBL 6cc/150 mg cartridges (Waters corp., Milford, MA, USA). Cartridges were conditioned with 12.5 mL of MeOH followed by drying under ~ 0.2 bar vacuum for 10 min. The cartridge was equilibrated with 7.5 mL MilliQ water before 125 mL of sample water (Laduviken or tap) was loaded and passed with a flow rate of approximately 2 mL/min. Immediately after the sample, the cartridge was washed with 5% v/v MeOH in MilliQ water, and then dried under ~ 0.2 bar vacuum for 10 min. The extract was eluted with 7.5 mL of MeOH, and then evaporated under a gentle air flow at 30°C. Both extracts were later reconstituted in 3 mL MilliQ water, and the water samples (SB1 – SB4) were spiked with abovementioned mixes (standard, unknown, and ILIS mix). Two different concentrations for unknowns were used. Final concentrations can be found in table S2 in SI.

## 2.3. Instrumental

Spiked water samples were analysed with Dionex UltiMate™ 3000 UHPLC system consisting of an RS Pump, RS Autosampler and RS Column Compartment (Thermo Fisher Scientific™, USA). The chromatographic separation was carried out on a Kinetex 2.6 µm, EVO C18, 100 Å, 150 × 3.0 mm column (Phenomenex®, Torrence, CA, USA). The mobile phases used was A: 0.1 % formic acid in MilliQ water with pH 2.7, and B: MeCN. Eluent gradient started with 5 % B, and then gradually increased to 100 % B over 20 minutes, was held at 100% B for 5 minutes and then decreased back to starting composition over 0.1 min. Between each injection, the column was equilibrated for 5 minutes. Sampler temperature was 15 °C, column oven temperature 40 °C, flow rate was 0.35 mL/min and the injection volume were 5 µL.

Samples were measured in positive ESI mode on a Q Exactive Orbitrap HRMS (Thermo Fisher Scientific™, USA), in two scan ranges;  $m/z$  60.0000 – 900.0000 and  $m/z$  100.0000 – 1 500.0000, with a mass resolution of 120 000. In the ion source, spray voltage was 3 500 V, capillary and probe heater temperature were 320 °C, max spray current was 100 µA and the S-lens RF level was 50 %. The gas parameters were sheath gas: 35, aux gas: 3, and spare gas: 0 (all in arbitrary units).

## 2.4. Data Treatment

Peaks from the chromatograms were integrated using the processing tool in XCalibur (Thermo Fisher Scientific™, USA), based on the  $m/z$  for each compound, with a mass tolerance of 10 ppm. From the processed files, information of the ions detected ( $[M+H]^+$ ,  $[M+Na]^+$ ,  $[M+NH_4]^+$ ,  $[M]^{2+}$ ) and their respective peak area and retention time were exported to Excel (Microsoft corp., USA) and saved as .csv files for further analysis in R (© The R Foundation).

The detected compounds were semi-quantified in R Studio, according to the principles described in section 1.5. If the predicted concentration ( $c_{predicted}$ ) was larger than the real concentration ( $c_{real}$ ), then

$$prediction\ error = \frac{c_{predicted}}{c_{real}}. \quad (5)$$

If instead the real concentration was larger than the predicted concentration, then

$$prediction\ error = \frac{c_{real}}{c_{predicted}}. \quad (6)$$

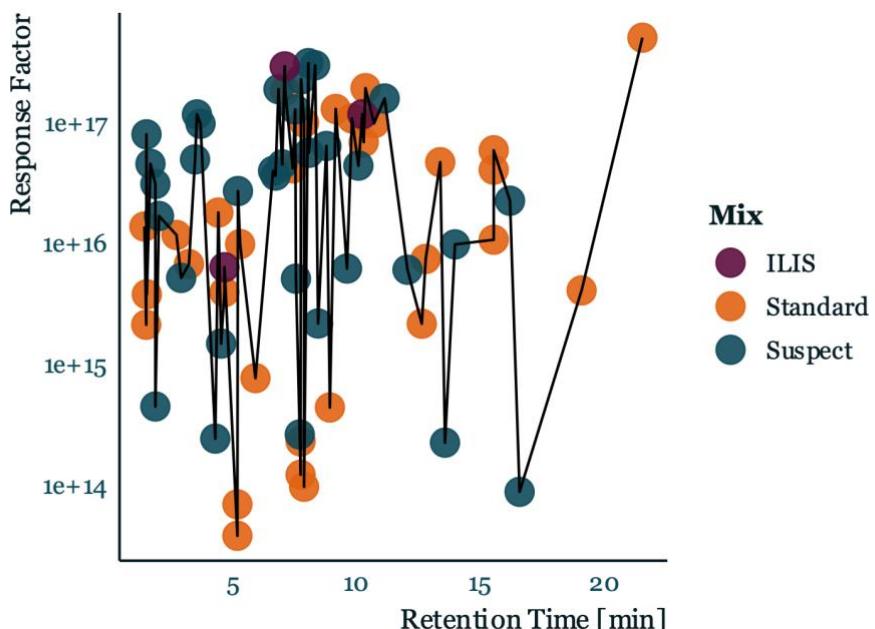
The code for the semi-quantification methods is available in supplementary information.

### 3. Results

A total of 67 compound for standards mix and suspect mix were selected from the NORMAN suspect list, to evaluate four semi-quantification strategies used in suspect screening. Selection was based partly on the chemicals' probability to end up in the surface water, and on their relevance in daily life, e.g. pharmaceuticals, pesticides, cosmetics, etc. There was also a desire to cover as wide chemical space as possible, especially in the aspect of chromatographic and ionisation properties, i.e. hydrophobicity and basicity. Thus, selection was also highly influenced by the compounds'  $\log P$  and  $pK_a$ . Due to technical reasons, two of the compounds, Cefoperazone and Avermectin B1a, were not available at the time for this thesis but will be included in the interlaboratory comparison. Before making the final decision of which compounds to include, test runs were performed at different concentrations to make sure that the substances were visible in the chromatogram. Three compounds, one from the standard mix and two from the suspect mix, were included even though they were not visible in the test runs. The standard compound, Chlormequat, is needed for one of the approaches in the interlaboratory trial. The two suspect compounds, Metolachlor-ESA and Metolachlor OXA was included as they were two of few TPs, that might be visible for other laboratories in the upcoming comparison. These three compounds are, however, excluded from the further discussion. Three additional compounds, that were seen in the test runs but not in the real runs, are also excluded from the discussion. These were Octocrylene, Butylamine and Melamine. The test runs also ensured that the concentrations chosen to spike the water samples were in the linear range. This is very important for the RF to be applicable in the semi-quantification methods used.

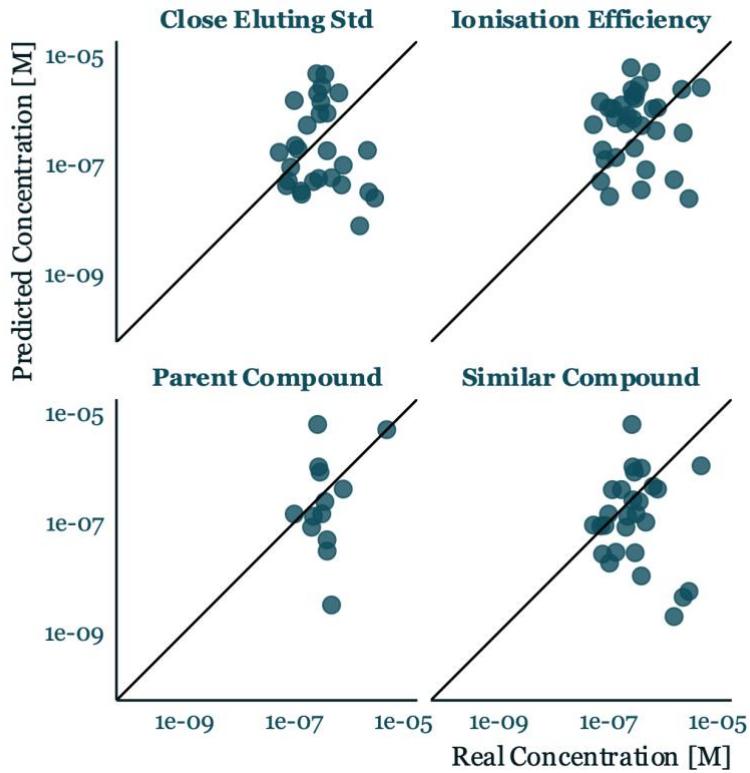
Additionally, three isotope labelled standards (ILIS mix) were added to the samples to adjust for any variation in injection volume. Spiked water samples (SB1 – SB4) were analysed and the 37 compounds in suspect mix were semi-quantified according to the proposed approaches.

As can be seen in figure 3, the standard as well as the suspect compounds are well distributed, with both retention time and response factor for both mixes dispersed over a wide range.

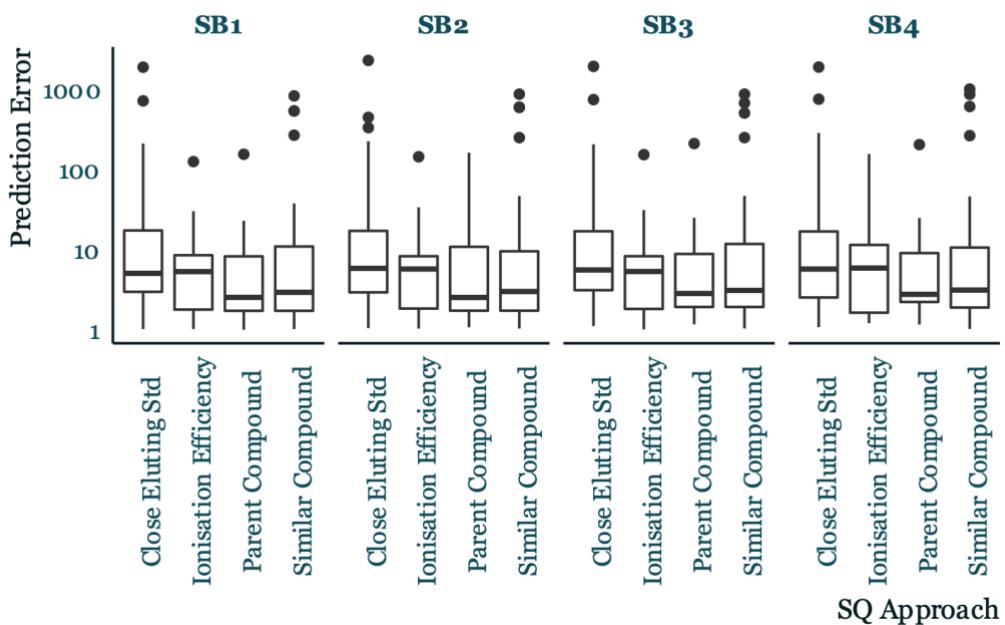


**Figure 3.** Plot displaying the distribution of suspect- and standard compounds over RT and RF for sample SB1. Distribution for the other samples did not deviate significantly.

Figure 4 displays how well the predicted concentrations correlate to the real concentrations, and in figure 5, a boxplot with the prediction errors for each approach is presented. The boxplot shows that the parent- and similar compound approaches give the lowest median error, while the close eluting standard approach appears to give the largest error. At the same time, all approaches give some unacceptable predictions, with over 1000-fold error for the close eluting standard and similar compound approach. The data in figure 4 are from one sample, however, data from the other three samples give similar correlation.



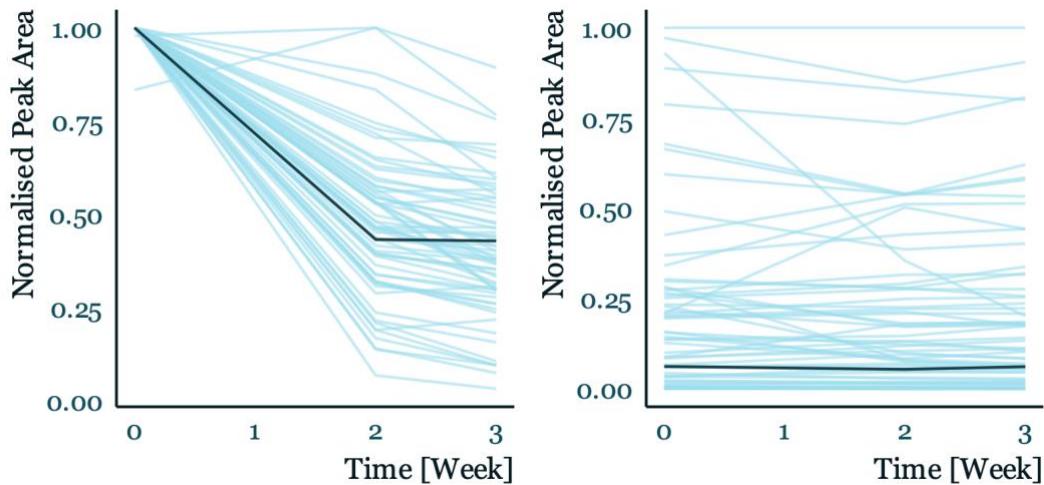
**Figure 4.** Plot of the predicted concentration vs the real concentration for each semi-quantification strategy used, for sample SB1. The data from the other samples were very similar to this data.



**Figure 5.** Boxplot displaying the prediction error for each approach and each sample.

Analysis of the samples was performed once a week for four weeks, to evaluate the stability of the compounds in water matrix. Stability testing is important for many reasons, in this case mainly to determine the suitability to send spiked water samples to the laboratories participating in the NORMAN collaborative trial. The stability is presented in figure 6, with two different approaches to normalise the peak areas.

Measurements from week 1 was ignored as all the intensities was much lower than the following weeks. Plots including this measurement can be found in SI, figure S1.



**Figure 6.** Stability plots showing how the peak area for each compound changes over time. To the left the measurement with the largest peak area for each compound is used to normalise that compounds peak areas for the other measurements. Almost all compounds that are used for the normalisation is from week zero, with two exceptions: Benzothiazole and 2-Methylbenzothiazole. This normalisation approach assumes that the instrument is stable over the whole analysis period. The right plot has instead normalised the peak areas relative to the compound with the largest peak area (Metformin), with the assumption that this substance is stable.

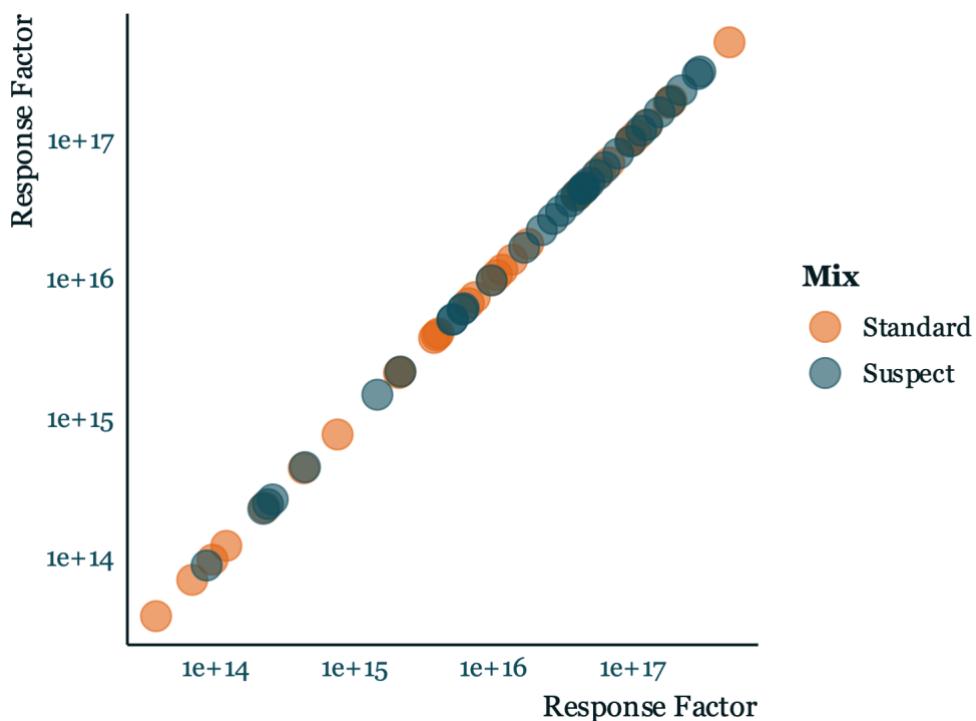
As seen, the different ways to normalise the peak areas give quite contradicting results. In the left-side plot in figure 6, almost all compounds appear to degrade quite rapidly, while in the right-side plot most compounds seem to be stable. The inconclusive results show that the stability needs to be further investigated during the interlaboratory trial with proper quantification approach.

## 4. Discussion

### 4.1. Compound Selection and Distribution

One of the main focuses for this study, and a very important aspect for the interlaboratory collaboration, was to select which analytes to include. As mentioned, the aim was to choose compounds with large variations in both chromatographic properties and relevance in the environment. The selection was also influenced by the semi-quantification methods that were to be used, i.e. parent – TP pairs were needed, as well as a good distribution between the standard and suspect compounds retention times. 18 of the standard compounds were already preselected by University of Athens; 2 of these were the ones that were not available at the time for the analysis. These 18 compounds are needed for calibration in one of the semi-quantification approaches suggested by their lab. Unfortunately, it was not possible to include that approach in this pre-study, as it has not yet been completely developed.

In figure 7, RF is plotted against RF for the standards and suspect analytes, to display the distribution of the two mixes. It can be seen that the majority of the analytes have quite high sensitivity. It would be desirable if more substances, especially from the suspect mix, would have had lower sensitivity. On the other hand, low responding compounds are generally harder to predict, and can also be harder to detect in the analysis, depending on the instrument used.



**Figure 7.** Plot of the RF vs RF for the suspect and standard compounds, showing the distribution of the substances with respect to their sensitivity.

Apart from wanting more analytes with lower response factor, the distribution between standards and suspect compounds seems to be quite satisfactory, based on the distribution of RF values. However, from figure 3, it can be deduced that some more late eluting suspect compounds would have been beneficial. After 12 min, 6 standards elute and 5 suspect compounds. The three data points at RT  $\approx$  15 min actually belong to the same compound, just different adducts.

There are no obvious compounds that could be interchanged between the suspect and standard mix. The two latest eluting substances seen in figure 3 (Ivermectin B1a and Nigericin)

both belong to the calibration mix from University of Athens and thus have to be in the standard mix. Similarly, the two lowest responding peaks in figure 7, are  $[M+Na]^+$  and  $[M+NH_4]^+$  adducts of Sucralose, so it would not be useful to substitute them for a suspect compound.

#### 4.2. Prediction Errors

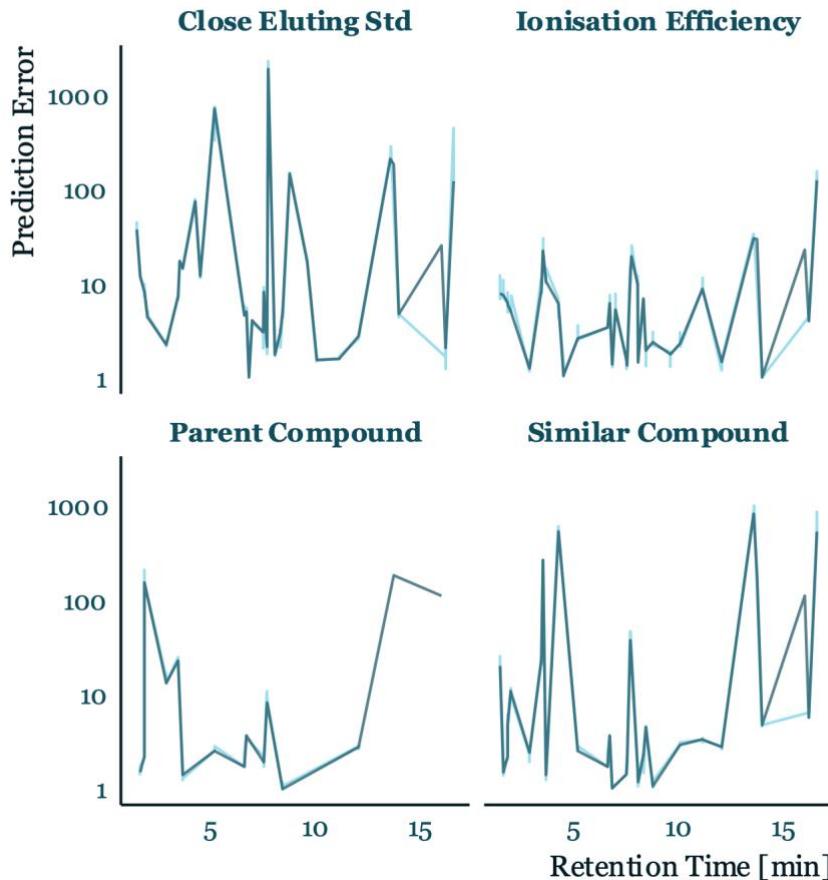
In figure 5, the boxplot displaying the prediction error for each semi-quantification approach, it would appear that best predictions are made with the parent compound approach. At least this approach is displaying the lowest median prediction errors across all the samples; however, this is not supported by the maximum errors or the mean errors, see table 1. By these criteria, the ionisation efficiency approach appears to be most suitable, even though the median error is the second worst of the four approaches. The closest eluting standard approach exhibits the poorest results, with up to over 2000-fold prediction error for Carbamazepine-10,11-epoxide. This is supported by both table 1 and figure 5.

**Table 1.** Mean-, median-, and max error for each sample and semi-quantification strategy.

Approach	Mean error	Median error	Max error	Sample
<b>Close Eluting Std</b>	101.6	7.4	1880.8	SB1
<b>Close Eluting Std</b>	113.5	6.8	2286.4	SB2
<b>Close Eluting Std</b>	101.7	7.4	1924.3	SB3
<b>Close Eluting Std</b>	104.0	7.4	1886.6	SB4
<b>Ionisation Efficiency</b>	11.1	6.1	125.3	SB1
<b>Ionisation Efficiency</b>	12.2	6.1	143.9	SB2
<b>Ionisation Efficiency</b>	11.8	5.5	153.0	SB3
<b>Ionisation Efficiency</b>	12.7	6.5	155.8	SB4
<b>Parent Compound</b>	34.3	2.8	184.2	SB1
<b>Parent Compound</b>	35.1	2.7	178.5	SB2
<b>Parent Compound</b>	36.1	3.0	210.2	SB3
<b>Parent Compound</b>	31.3	2.8	203.9	SB4
<b>Similar Compound</b>	84.1	3.4	825.2	SB1
<b>Similar Compound</b>	89.2	3.1	867.3	SB2
<b>Similar Compound</b>	87.6	3.3	868.1	SB3
<b>Similar Compound</b>	100.1	3.3	1004.1	SB4

#### 4.2.1. Correlation of the Errors with Chromatographic Properties

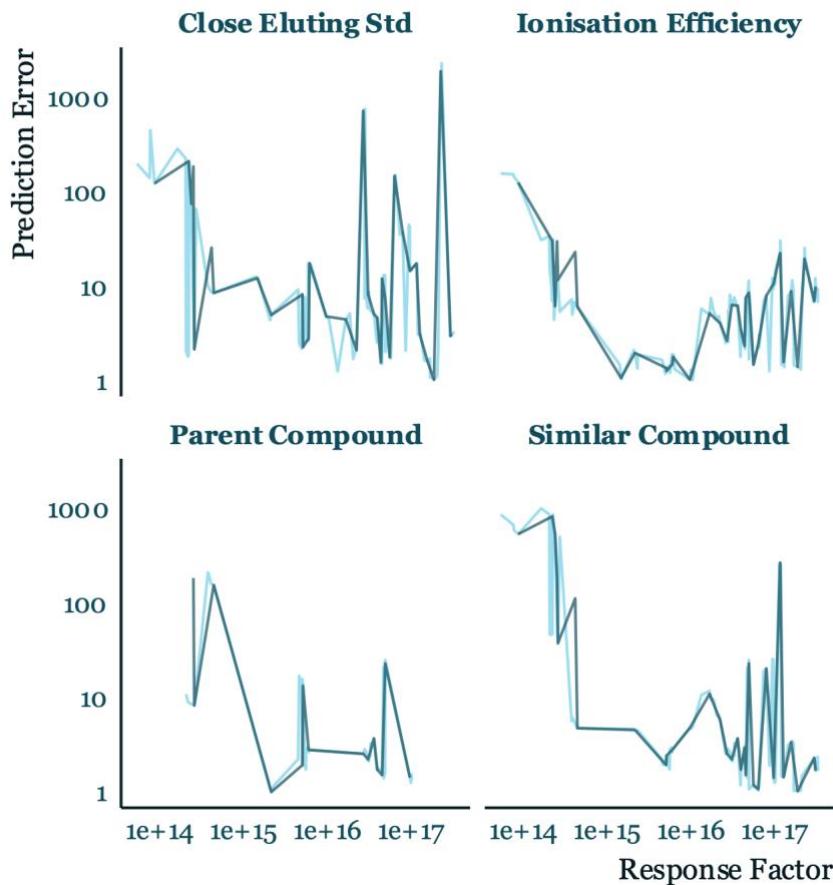
Analysis of the data revealed that there is insignificant correlation between the retention time and the prediction error, as seen in figure 8. The RT for the analytes with the highest errors ranged from 1.85 min for Atrazine-desethyl-desisopropyl-2-hydroxy to 16.58 min for Chlorporifos.



**Figure 8.** Plot showing that there is no trend to be seen with regards to the prediction error and the retention time. As seen, the pattern is basically the same across all samples.

All of the compounds yielding largest prediction errors are quite weak bases, with  $pK_a$  values from -4.2 to 4.6. This weak basicity makes the compounds harder to ionise, thus they yield a lower response factor, which in turn makes them harder to predict. This has been shown previously by both Liigand et al. (2020),<sup>23</sup> and by Kruve & Kaupmees (2017) for the ionisation efficiency prediction models.<sup>28</sup> Though this connection has been shown for *IE* prediction models only, it appears to be the case for the other approaches too, as seen in figure 9. The trend is therefore most distinct in the *IE* approach, but is also seen in the other approaches. However, some of the substances with unacceptably high error still has high response factor, which is especially clear in the closest eluting standard and the similar compound approaches. In fact, three out of the five substances with highest error in the closest eluting standard approach are in the high RF range. These are, ordered based on RF, Atrazine-desisopropyl ( $\text{mean RF} \approx 3e+16$ ,  $\text{mean error} \approx 630\text{-fold}$ ), Simazine ( $\text{mean RF} \approx 6e+16$ ,  $\text{mean error} \approx 145\text{-fold}$ ) and Carbamazepine-10,11-epoxide ( $\text{mean RF} \approx 2e+17$ ,  $\text{mean error} \approx 2000\text{-fold}$ ). In the similar compound approach, only one of the top five compounds with highest error also had high sensitivity, namely Omethoate ( $\text{mean RF} \approx 1e+17$ ,  $\text{mean error} \approx 260\text{-fold}$ ). The trend that the error decreases as the response factor increases, seems correct up to a certain value of RF, where the error starts increasing again. This is seen for all strategies, including the ionisation

efficiency, and thereby indicates that the prediction error cannot be explained by response factor alone.



**Figure 9.** Plot displaying an overall trend that analytes with low sensitivity gives high prediction errors. This trend is similar for all four samples and approaches.

#### 4.2.2. Adduct Formation and Semi-Quantification Error

Some of the substances with the highest errors formed adducts, e.g. Atrazine-2-hydroxy formed sodium adduct, Efavirenz, Chlorothiazide and 2-Hydroxybenzothiazole formed ammonium adducts, and Carbamazepine-10,11-epoxide and Omethoate formed both sodium and ammonium adducts. Of these, only the Chlorothiazide ammonium adduct had higher intensity than the protonated ion, which was the peak that was integrated and used in the semi-quantification. Therefore, this might explain high prediction error for Chlorothiazide. The Efavirenz ammonium ion was around 20% of the intensity of  $[M+H]^+$  ion, whereas the adducts of the other analytes were less than 10%. Thus, they were deemed as not significant to explain the high prediction errors. The mass spectra of the compounds with highest error was also investigated, to determine whether any insource fragmentation occurred. Though there were some compounds that displayed fragmentation pattern, none of them were significant, i.e. less than 10 % intensity.

### 4.3. Structural Similarity

In table S4 in supplementary information, the suspect compounds are listed, together with the closest eluting/parent/most similar standard compound. In table S5, structures of all compounds that were visible in the analysis is displayed. With the exception of the parent compound approach, many of the compounds with the highest error across the semi-quantification methods are the same. E.g. Efavirenz and Chlorpyrifos, with prediction error

ranging from approximately 30- to 100-fold and 120- to 670-fold, respectively, are among the five highest substances with highest error for all approaches except the parent compound one. It should be mentioned that another compound with the same thiophosphate group as in Chlorpyrifos, previously has shown similar poor prediction results. In unpublished work by Kiefer et al.,<sup>29</sup> a similar compound, Diazinon, was the one with the highest error for both parent compound approach and ionisation efficiency predictions. This leads to suspicions that something with the thiophosphate moiety makes molecules harder to semi-quantify potentially referring to uninvestigated aspects in the ionisation mechanism. It might also explain why Omethoate is among the compounds with highest error, as it possess a similar phosphate moiety.

A qualitative analysis was performed for the five substances giving the highest error, to evaluate their respective similarity. For the parent compound approach, this analysis was done for three compounds only, since there were only three that gave error that were deemed unacceptable, i.e. more than 10-fold. In all, the closest eluting standard, *IE*, and similar structure approaches had 12, 6-11, and 7-8 compounds respectively that gave error higher than acceptable.

#### 4.3.1. Parent Compound Approach

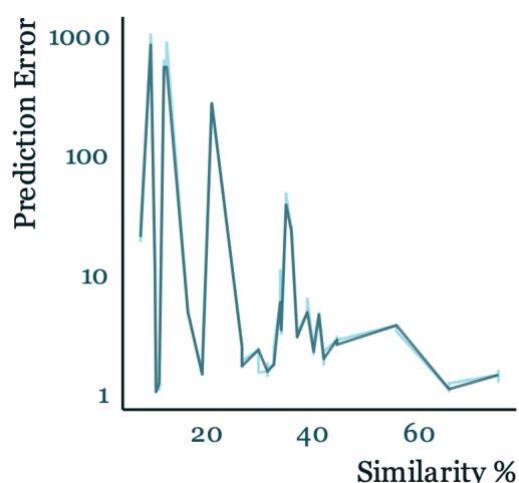
For the parent compound approach, the TPs look fairly similar to their parent compound, as seen in table 2. Yet, according to the similarity online tool described in section 1.5.1.,<sup>24</sup> the Atrazine TPs both shared more similarity with Guanylurea than with Atrazine. Qualitatively, significant differences in the functional groups can be observed as well for the TPs in table 1. Both the isopropyl- and the ethyl group has been cleaved in Atrazine-desethyl-desisopropyl as well as for Atrazine-desethyl-desisopropyl-2-hydroxy. For the latter TP, the chlorine atom has also been exchanged for a hydroxy group. Similarly, as TCMTB degrades to 2-Aminobenzothiazole, it loses quite a large moiety, which might explain why this TP gives such a high prediction error.

Comparing the hydrophobicity, it is seen that it decreases for all of the TPs compared to the parent compounds. However, this is more pronounced for the Atrazine TPs than for the transformation product of TCMTB. From the hydrophobicity aspect, Atrazine-desethyl-desisopropyl-2-hydroxy indeed seems to be more similar to Guanylurea than Atrazine. On the contrary, comparison of the basicity shows quite the opposite. Even though the  $pK_a$  values

differs between parent and TP, they are all in the range of weak bases.

#### 4.3.2. Similar Compound Approach

As seen in table 3, none of the most similar structures correlates well with the structure it is supposedly most similar to, at least not for the five analytes with the highest prediction error. Basically, the suspect compound shares some aromaticity with its most similar standard compound, and some common atoms can be found in the structures. The similarity score is ranging from 9.5 % (for Efavirenz and Imazalil) to 35 % (for 2-Hydroxybenzothiazole and Benzotriazole), which is not such a high score. The lower the similarity score is, the higher the maximum expected error becomes. Still, none of the suspects

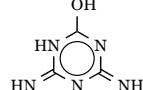
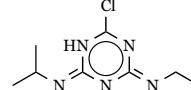
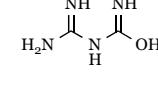
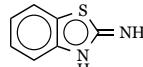
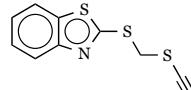
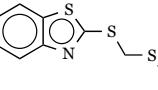
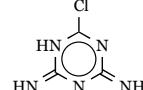
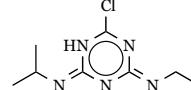
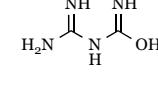


**Figure 10.** Plot showing the correlation between similarity score and prediction error. The trend appears to be that the prediction error decreases as the similarity score increases.

are within the maximum expected error range. In fact, the mean error for Efavirenz is approximately 50x higher than the maximum expected error, and for 2-Hydroxybenzothiazole the mean error is approximately 5x higher than expected.

Based on the results in this study, it seems to be a correlation between the similarity score and the error, as seen in figure 10, which is hardly surprising. However, there are exceptions, e.g. Metformin, which is 7.6 % similar to its most similar compound Caffeine. Still, the mean prediction error is approximately 20-fold, i.e. much less than for Efavirenz although the similarity score is lower. This is likely to result from probability in obtaining similar response factors for two well ionising compounds, not from significant structural similarities.

**Table 2.** Structural, hydrophobicity and basicity comparison of the transformation product and the parent compound, with average error over 10-fold. The most similar compound is also included to make the discussion easier to follow.  $pK_a$ - and  $\log P$  values were all calculated in ChemAxon.

TP	Parent	Most similar std
Atrazine-desethyl-desisopropyl-2-hydroxy   Mean error: 182.6 $\log P$ : -4.5 $pK_a$ : 3.1	Atrazine  	Guanylurea  
2-Aminobenzothiazole   Mean error: 23.4 $\log P$ : 1.9 $pK_a$ : 4.48	TCMTB  	TCMTB  
Atrazine-desethyl-desisopropyl   Mean error: 14.8 $\log P$ : -0.25 $pK_a$ : 4.58	Atrazine  	Guanylurea  

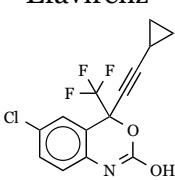
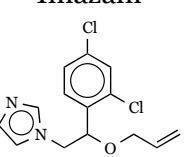
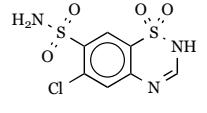
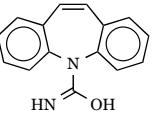
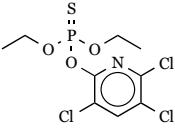
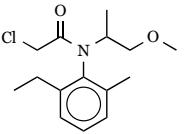
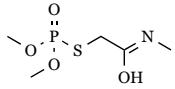
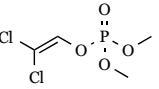
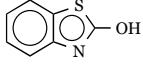
#### 4.3.3. Closest Eluting Standard Approach

The similarity between the suspect compound and the closest eluting standard for this approach is the worst across all the approaches. In table 4, it is seen that there is practically no similarity at all between the pairs. However, this strategy is not based on finding compounds with structural similarities, but rather to find compound with similar retention time.

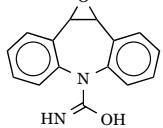
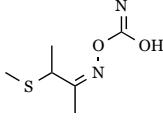
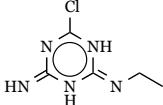
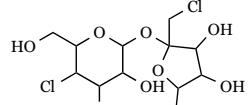
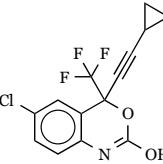
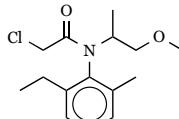
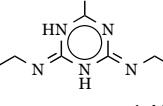
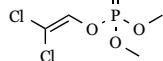
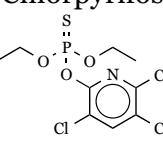
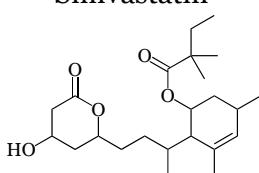
All of the retention times of the standard compounds are close to the suspects, from a couple of seconds (Carbamazepine-10,11-epoxide and Atrazine-desisopropyl), to 12 s for Efavirenz. The exception is Chlorpyrifos, which is eluting more than one minute after its closest neighbour Simvastatin. As previously discussed, more late eluting compounds in the mix would have been desirable. If this had been implemented, there would probably have been a standard compound eluting closer in time to the suspect. However, it does not appear as the closer RT between suspect and standard give lower errors, but rather the opposite. This approach is overall the worst of the ones tested in this study, as can be seen concluded by

looking at the high errors. Most probably, the response factor which is used for the semi-quantification depends on many factors that are not accounted for by this approach.

**Table 3.** Structural, hydrophobicity and basicity comparison of the suspect compound and its most similar standard. The mean error of the suspect is also compared with the maximum expected error.  $\log P$  and  $pK_a$  were calculated in ChemAxon, similarity score and maximum expected error were from <http://dsfp.chem.uoa.gr/semi quantification/>.

Suspect compound	Most similar standard compound
<b>Efavirenz</b>  Mean error: 891.2 $\log P$ : 4.5 $pK_a$ : -1.5	<b>Imazalil</b>  Similarity score: 9.5 % Maximum expected error: 18.6 $\log P$ : 3.8 $pK_a$ : 6.5
<b>Chlorothiazide</b>  Mean error: 562.0 $\log P$ : -0.44 $pK_a$ : 1.2	<b>Carbamazepine</b>  Similarity score: 12.0 % Maximum expected error: 15.1 $\log P$ : 2.8 $pK_a$ : -3.8
<b>Chlorpyrifos</b>  Mean error: 666.3 $\log P$ : 4.8 $pK_a$ : -4.2	<b>Metolachlor</b>  Similarity score: 12.5 % Maximum expected error: 15.1 $\log P$ : 3.5 $pK_a$ : -4.1
<b>Omethoate</b>  Mean error: 257.5 $\log P$ : -0.55 $pK_a$ : NA	<b>Dichlorvos</b>  Similarity score: 21.0 % Maximum expected error: 11.9 $\log P$ : 1.4 $pK_a$ : NA
<b>2-Hydroxybenzothiazole</b>  Mean error: 44.4 $\log P$ : 2.5 $pK_a$ : -1.3	<b>Benzotriazole</b>  Similarity score: 35.0 % Maximum expected error: 9.2 $\log P$ : 1.3 $pK_a$ : 0.22

**Table 4.** Comparison of structure as well as RT for the suspect compound and the closest eluting standard.  
 \*Expert pKa estimation for Carbamazepine-10,11-epoxide, Dichlorvos and Simvastatin was 1 to 4 for Carbamazepine-10,11-epoxide, around -7 for Simvastatin and approximately 10 units lower for Dichlorvos.

Suspect compound	Closest eluting standard
<b>Carbamazepine-10,11-epoxide</b>  Mean error: 1994.5 RT: 7.74 min $\log P$ : 2.0 $pK_a$ : NA* Mean RF: $2.3e+17$	<b>Butocarboxim</b>  RT: 7.72 min $\log P$ : 1.3 $pK_a$ : 1.33 Mean RF: $1.1e+14$
<b>Atrazine-desisopropyl</b>  Mean error: 635.0 RT: 5.19 min $\log P$ : 0.39 $pK_a$ : 4.4 Mean RF: $2.7e+16$	<b>Sucralose [M+Na]<sup>+</sup></b>  RT: 5.17 min $\log P$ : -0.47 $pK_a$ : - Mean RF: $4.4e+13$
<b>Efavirenz</b>  Mean error: 231.5 RT: 13.57 min $\log P$ : 4.5 $pK_a$ : -1.5 Mean RF: $2.0e+14$	<b>Metolachlor</b>  RT: 13.37 min $\log P$ : 3.5 $pK_a$ : -4.1 Mean RF: $4.7e+16$
<b>Simazine</b>  Mean error: 145.2 RT: 8.77 min $\log P$ : 1.8 $pK_a$ : 4.2 Mean RF: $6.4e+16$	<b>Dichlorvos</b>  RT: 8.91 min $\log P$ : 1.4 $pK_a$ : NA* Mean RF: $4.5e+14$
<b>Chlorpyrifos</b>  Mean error: 226.5 RT: 16.58 min $\log P$ : 4.8 $pK_a$ : -4.2 Mean RF: $7.3e+13$	<b>Simvastatin</b>  RT: 15.54 min $\log P$ : 4.5 $pK_a$ : NA* Mean RF: $2.0e+16$

#### 4.3.4. Ionisation Efficiency Approach

This approach did not use the response factor of any specific standard compound to quantify the suspect, but rather predicted the ionisation efficiency. The prediction was based on 2D molecular descriptors, similar as to a fingerprint of the molecule. These descriptors aim to give similar results as 3D optimisation, but at a realistic time scale. As the model was developed, it was noticed that structural parameters were of most importance.<sup>23</sup> Among the 15 key features was the number of hydrogen bonds and the number of nitrogen atoms in the molecule. These can be related to hydrophobicity and basicity, respectively,<sup>23</sup> which are the two main properties discussed in this thesis. As seen in table 5, logP and pK<sub>a</sub> of the compounds with the highest errors are displayed. Theoretically, higher logP values should yield lower prediction error; however, this does not seem to be the case for these substances. This might be because the relatively high logP values are overshadowed by the very weak basicity exhibited by the compounds.

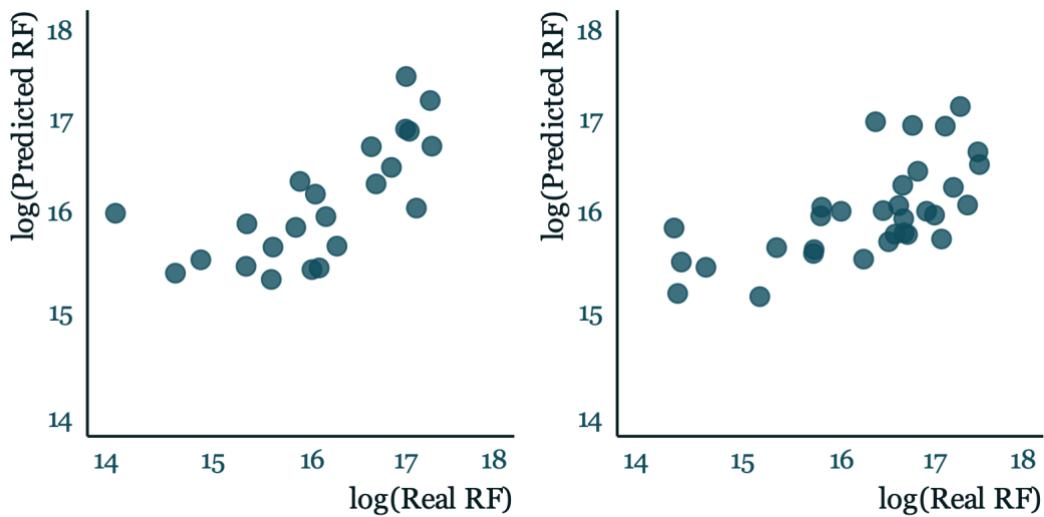
**Table 5.** The substances with the highest error for the ionisation efficiency approach, with each compounds respective RF, predicted RF, logP, pK<sub>a</sub> and mean error. logP and pK<sub>a</sub> was calculated in ChemAxon.

\*Expert pK<sub>a</sub> estimation for Omethoate and Carbamazepine-10,11-epoxide, was between 1 to 4 for Carbamazepine-10,11-epoxide and approximately -2 for Omethoate.

Suspect compound	log (Real RF)	log (Predicted RF)	logP	pK <sub>a</sub>	Mean error
Chlorpyrifos	13.9	16.0	4.8	-4.2	144.5
Efavirenz	14.3	15.8	4.5	-1.5	31.5
Omethoate	17.1	15.7	-0.55	NA*	24.0
Carbamazepine-10,11-epoxide	17.4	16.8	2.0	NA*	20.8
2-Hydroxy benzothiazole	14.4	15.5	2.5	-1.3	12.7

In figure 11, the logarithm of the predicted response factor is plotted against the logarithm of the real response factor. As seen, the data points all follow some trendline, albeit the correlation is somewhat better for the standard substances, i.e. the left-hand plot. Obviously,

the error is increasing with the difference between the predicted and real RF, as seen in table 5.



**Figure 11.** Plots showing the predicted response factor vs the real response factor. The left is for the standard compounds and the right plot is for the suspect compounds.

#### 4.4. Stability

The results from the stability testing gave somewhat inconclusive results, as seen when comparing the two normalisation approaches in figure 6. The different ways to normalise the peak areas give quite opposite results, and neither of the assumptions that were made were necessarily correct. Actually, the measurements from week 1 would suggest that the instrument in fact was not stable, and thus, the other normalisation approach would appear to be better, i.e. the right-hand plot. However, the stability of Metformin in water is unknown, and thus the normalisation against Metformin might not be valid. Therefore, the stability needs additional investigation and quantification based on analytical standards.

## 5. Conclusion

Development of new analytical techniques is a prerequisite for increasing environmental awareness. For this purpose, non-targeted screening of environmental samples is being employed more and more. This is especially true in the search for emerging contaminants in water samples. Though semi-quantification methods for non-targeted analysis are available, these are not yet standardised. A desire to validate semi-quantitative non-targeted analysis led the NORMAN Network to arrange an interlaboratory comparison on the subject, organised by Stockholm University and University of Athens. As part of this thesis, four semi-quantification approaches proposed for collaborative trial were tested.

The results show quite clearly that the closest eluting standard approach is the worst performing strategy, no matter which criteria it is based on. This approach gave both the highest mean and median error, with up to 2000-fold maximum error. It also gave unacceptable high errors (over 10-fold) for the largest number of compounds. This is quite unfortunate, since this strategy was the only one of the four that could be easily used in true non-targeted screening as it does not require a tentative structure.

Comparing the three other strategies is not as straight forward. The similar compound and parent compound approaches present the lowest median error, with the parent approach having slightly lower median error. The parent compound strategy was also able to give satisfying predictions for all compounds except three. However, this approach does not work for all compound, as it is applicable to degradation products. The similar compound method is more universal, as most substances are to some extent similar to another substance. Still, the maximum errors for this approach is up to 1000-fold, which thereby gives quite high mean errors. Across the samples, this approach had undesirable results for 7-8 compounds. The ionisation efficiency method on the other hand exhibited the lowest mean error of all approaches, and the second lowest number of compounds with over 10-fold error. However, the median error was the second worst for this approach. Though the *IE* approach seems to be the most promising method, the fact that it only works for protonated or deprotonated molecules is unfortunate, as it means that it will not work for all compounds.

In conclusion, no one approach appears to be sufficient as a universal semi-quantification approach, but rather different strategies can to be used in parallel where applicable.

## **6. Future Challenges**

The results presented in this thesis need to, and will be, further investigated. In the upcoming NORMAN collaborative trial, the results from different laboratories will be compared, and two additional semi-quantification strategies will be tested. More data input will, hopefully, show clearer trends, and thus facilitate evaluation of the approaches.

For future work, further development of ionisation efficiency prediction models is encouraged. Especially, further investigation in predicting the low sensitivity substances needs to improve prediction accuracy. Additionally, adding the prediction possibility for sodium or ammonium adducts is beneficial to widen the application range method.

## **7. Acknowledgements**

I would like to say special thanks to my supervisor Anneli Kruve, first of all for giving me the amazing opportunity to be a part of the NORMAN collaborative trial. Secondly, for always being willing to help me and explain when I was in doubt. Thirdly, for all the ice cream along this project.

I would also like to thank Claudia Möckel and Merle Plassman for all technical support and help with the instruments. Further, I thank Miklós Mohai, as his previous work helped with the compound selection, Josefina Carlsson and Ulrika Nilsson from the textile lab group for letting me raid their chemical cupboard, and Dr. Karl Kaupmees for his expert estimations of  $pK_a$  values.

Special thanks also to Emma Palm, Riccardo Costalunga, Amina Souihé and Anselm Okolo. Finally, I thank Jonathan Benskin for taking the time to be my opponent.

## 8. References

- (1) NORMAN Network. Suspect list <https://www.normandata.eu/normansusdat/>.
- (2) Pieke, E. N.; Granby, K.; Trier, X.; Smedsgaard, J. A Framework to Estimate Concentrations of Potentially Unknown Substances by Semi-Quantification in Liquid Chromatography Electrospray Ionization Mass Spectrometry. *Analytica Chimica Acta* **2017**, *975*, 30–41. <https://doi.org/10.1016/j.aca.2017.03.054>.
- (3) Kruve, A. Strategies for Drawing Quantitative Conclusions from Nontargeted Liquid Chromatography–High-Resolution Mass Spectrometry Analysis. *Anal. Chem.* **2020**, *92* (7), 4691–4699. <https://doi.org/10.1021/acs.analchem.9b03481>.
- (4) World Health Organisation. Drinking Water <https://www.who.int/en/news-room/fact-sheets/detail/drinking-water>.
- (5) Bukola, D.; Zaid, A. Consequences of Anthropogenic Activities on Fish and the Aquatic Environment. *Poult Fish Wildl Sci* **2015**, *03* (02). <https://doi.org/10.4172/2375-446X.1000138>.
- (6) Schwarzenbach, R. P. The Challenge of Micropollutants in Aquatic Systems. *Science* **2006**, *313* (5790), 1072–1077. <https://doi.org/10.1126/science.1127291>.
- (7) Richardson, S. D.; Kimura, S. Y. Water Analysis: Emerging Contaminants and Current Issues. *Anal. Chem.* **2020**, *92* (1), 473–505. <https://doi.org/10.1021/acs.analchem.9b05269>.
- (8) Kruve, A. Semi-quantitative Non-target Analysis of Water with Liquid Chromatography/High-resolution Mass Spectrometry: How Far Are We? *Rapid Commun Mass Spectrom* **2019**, *33* (S3), 54–63. <https://doi.org/10.1002/rcm.8208>.
- (9) Bletsou, A. A.; Jeon, J.; Hollender, J.; Archontaki, E.; Thomaidis, N. S. Targeted and Non-Targeted Liquid Chromatography-Mass Spectrometric Workflows for Identification of Transformation Products of Emerging Pollutants in the Aquatic Environment. *TrAC Trends in Analytical Chemistry* **2015**, *66*, 32–44. <https://doi.org/10.1016/j.trac.2014.11.009>.
- (10) Kiefer, K.; Müller, A.; Singer, H.; Hollender, J. New Relevant Pesticide Transformation Products in Groundwater Detected Using Target and Suspect Screening for Agricultural and Urban Micropollutants with LC-HRMS. *Water Research* **2019**, *165*, 114972. <https://doi.org/10.1016/j.watres.2019.114972>.
- (11) European Commision. Directive of the European Parliament and of the Council: Amending Directives 2000/60/EC and 2008/105/EC as Regards Priority Substances in the Field of Water Policy. January 31, 2012.
- (12) United States Environmental Protection Agency. Toxic and priority pollutants under the clean water act <https://www.epa.gov/eg/toxic-and-priority-pollutants-under-clean-water-act#priority>.
- (13) Hollender, J.; Schymanski, E. L.; Singer, H. P.; Ferguson, P. L. Nontarget Screening with High Resolution Mass Spectrometry in the Environment: Ready to Go? *Environ. Sci. Technol.* **2017**, *51* (20), 11505–11512. <https://doi.org/10.1021/acs.est.7b02184>.
- (14) Schymanski, E. L.; Jeon, J.; Gulde, R.; Fenner, K.; Ruff, M.; Singer, H. P.; Hollender, J. Identifying Small Molecules via High Resolution Mass Spectrometry: Communicating Confidence. *Environ. Sci. Technol.* **2014**, *48* (4), 2097–2098. <https://doi.org/10.1021/es5002105>.
- (15) Hug, C.; Ulrich, N.; Schulze, T.; Brack, W.; Krauss, M. Identification of Novel Micropollutants in Wastewater by a Combination of Suspect and Nontarget Screening. *Environmental Pollution* **2014**, *184*, 25–32. <https://doi.org/10.1016/j.envpol.2013.07.048>.
- (16) Schymanski, E. L.; Singer, H. P.; Slobodnik, J.; Ipolyi, I. M.; Oswald, P.; Krauss, M.; Schulze, T.; Haglund, P.; Letzel, T.; Grosse, S.; Thomaidis, N. S.; Bletsou, A.; Zwiener, C.; Ibáñez, M.; Portolés, T.; de Boer, R.; Reid, M. J.; Onghena, M.; Kunkel, U.; Schulz, W.; Guillou, A.; Noyon, N.; Leroy, G.; Bados, P.; Bogialli, S.; Stipaničev, D.; Rostkowski, P.; Hollender, J. Non-Target Screening with High-Resolution Mass Spectrometry: Critical

- Review Using a Collaborative Trial on Water Analysis. *Anal Bioanal Chem* **2015**, *407* (21), 6237–6255. <https://doi.org/10.1007/s00216-015-8681-7>.
- (17) Waters. Oasis Solid-Phase Extraction Products. Waters April 2020.
- (18) NORMAN Network of reference laboratories, research centres and related organisations for monitoring of emerging environmental substances <https://www.norman-network.com/>.
- (19) NORMAN Network. NORMAN Joint Programme of Activities (JPA 2020). February 2020.
- (20) Guthrie, W. F. Interlaboratory Comparisons. *Wiley StatRef: Statistics Reference Online* 7. <https://doi.org/10.1002/9781118445112.stat04148>.
- (21) Vangel, M. G. Interlaboratory Studies. *Wiley StatRef: Statistics Reference Online* **2014**, 6. <https://doi.org/10.1002/9781118445112.stat07619>.
- (22) Kruve, A.; Aalizadeh, R.; Alygizakis, N.; Thomaidis, N. S.; Malm, L. Interlaboratory Comparison on Strategies for Semi-Quantitative Non-Targeted LC-ESI-HRMS <https://www.norman-network.net/sites/default/files/files/QA-QC%20Issues/Invitation%20letter%20JPA%202020%20semi-quant%20inter%20lab%20%28002%29.pdf>.
- (23) Liigand, J.; Wang, T.; Kellogg, J.; Smedsgaard, J.; Cech, N.; Kruve, A. Quantification for Non-Targeted LC/MS Screening without Standard Substances. *Sci Rep* **2020**, *10* (1), 5808. <https://doi.org/10.1038/s41598-020-62573-z>.
- (24) Similar compound finder <http://dsfp.chem.uga.edu/semitquantification/>.
- (25) Dahal, U. P.; Jones, J. P.; Davis, J. A.; Rock, D. A. Small Molecule Quantification by Liquid Chromatography-Mass Spectrometry for Metabolites of Drugs and Drug Candidates. *Drug Metab Dispos* **2011**, *39* (12), 2355–2360. <https://doi.org/10.1124/dmd.111.040865>.
- (26) Quantem Analytics <https://app.quantem.co/>.
- (27) Rousis, N. I.; Gracia-Lor, E.; Zuccato, E.; Bade, R.; Baz-Lomba, J. A.; Castrignanò, E.; Causanilles, A.; Covaci, A.; de Voogt, P.; Hernández, F.; Kasprzyk-Hordern, B.; Kinyua, J.; McCall, A.-K.; Plósz, B. Gy.; Ramin, P.; Ryu, Y.; Thomas, K. V.; van Nuijs, A.; Yang, Z.; Castiglion, S. Wastewater-Based Epidemiology to Assess Pan-European Pesticide Exposure. *Water Research* **2017**, *121*, 270–279. <https://doi.org/10.1016/j.watres.2017.05.044>.
- (28) Kruve, A.; Kaupmees, K. Predicting ESI/MS Signal Change for Anions in Different Solvents. *Anal Chem* **2017**, *89* (9), 5079–5086. <https://doi.org/10.1021/acs.analchem.7b00595>.
- (29) Kiefer, K.; Signer, H.; Hollender, J.; Kruve, A. Benchmarking of the Semi-Quantification Strategies for the Non-Targeted Screening of Micropollutants and Transformation Products in Groundwater. Unpublished.

## 9. Supplementary Information

**Table S1.** All the compounds with their corresponding solvent, masses and concentrations. Note that for nigericin, 10,11-dihydro-10-hydroxycarbamazepine, carbamazepine-10,11-epoxide and caffeine-<sup>13</sup>C<sub>3</sub> no mass for the solvent is available as they were purchased as solutions with concentration of 1mg/mL. Atrazine-desethyl-desisopropyl-2-hydroxy were prepared in a 100 mL volumetric flask and was not weighed.

Compound	<i>m</i> <sub>compound</sub> [mg]	Solvent	<i>m</i> <sub>solvent</sub> [mg]	<i>c</i> [M]
<b>Guanylurea</b>	2.4	H <sub>2</sub> O	10007.5	0.00158
<b>Amitrole</b>	8.8	MeCN	7789.6	0.01056
<b>Histamine</b>	11.1	MeCN	7786.2	0.01008
<b>Chlormequat</b>	10.7	MeCN	7786.2	0.00683
<b>Methamidophos</b>	17.5	MeCN	7763.4	0.01255
<b>Vancomycin</b>	1.7	H <sub>2</sub> O	9975.4	0.00012
<b>Trichlorfon</b>	8	MeCN	7751	0.00315
<b>Butocarboxim</b>	10.4	MeCN	7761.6	0.00554
<b>Dichlorvos</b>	14.7	MeCN	7718.2	0.00677
<b>Tylosin</b>	10.9	MeCN	7785.7	0.00108
<b>Rifaximin</b>	9.7	MeCN	7783.6	0.00125
<b>Spinosyn A</b>	10.2	MeCN	7810.8	0.00140
<b>Emamectin B1a</b>	12.1	MeCN	7778.4	0.00121
<b>Nigericin</b>	5	DMSO:EtOH 1:1	-	0.00670
<b>Ivermectin B1a</b>	11.3	MeCN	7733.8	0.00131
<b>TCMTB</b>	10.4	MeCN	7791.6	0.00440
<b>Atrazine</b>	8.8	MeCN	7796.8	0.00411
<b>Octocrylene</b>	19.7	MeCN	7776	0.00551
<b>Clarithromycin</b>	10.5	MeCN	7763.2	0.00142
<b>Aspartame</b>	13.4	H <sub>2</sub> O	9980.4	0.00455
<b>Simvastatin</b>	9.9	MeCN	7782.6	0.00239
<b>Sucralose</b>	13.7	MeCN	7766.8	0.00349
<b>Irgarol</b>	7	MeCN	7777	0.00279
<b>Caffeine</b>	13.4	MeCN	7771.5	0.00698
<b>Carbamazepine</b>	10.9	MeCN	7750.7	0.00371
<b>Benzotriazole</b>	14.6	MeCN	7707.3	0.01250
<b>Metolachlor</b>	12	MeCN	7759.2	0.00428
<b>Imazalil</b>	9.9	MeCN	7757.1	0.00338
<b>Atrazine- desethyl</b>	11.1	MeCN	7795.8	0.00596
<b>Atrazine-desethyl-2-hydroxy</b>	2.2	H <sub>2</sub> O:0.1% FA 5:1	11957.8	0.00108
<b>Atrazine-desisopropyl</b>	9.2	MeCN	7766.3	0.00536
<b>Atrazine-desethyl-desisopropyl</b>	2.3	DMSO	11009.6	0.00158

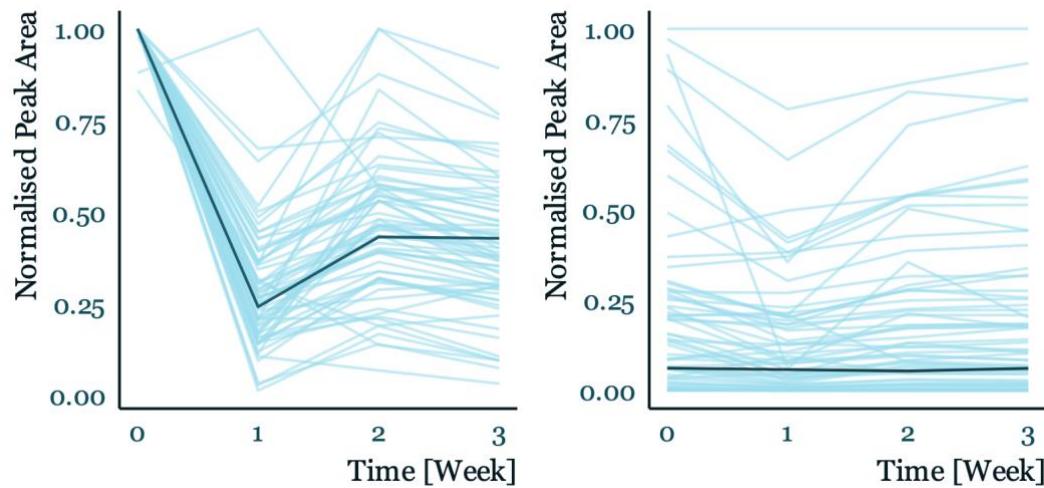
Compound	$m_{\text{compound}}$ [mg]	Solvent	$m_{\text{solvent}}$ [mg]	$c$ [M]
<b>Atrazine-desethyl-</b> <b>desisopropyl-2-</b> <b>hydroxy</b>	2.4	H <sub>2</sub> O:0.1% FA 1:10	-	0.00019
<b>Atrazine-desisopropyl-</b> <b>2-hydroxy</b>	2.3	H <sub>2</sub> O:0.1% FA 5:1	12028.8	0.00123
<b>2-(Methylthio)</b> <b>benzothiazole</b>	12.5	MeCN	7764.2	0.00848
<b>Progesterone</b>	11.8	MeCN	7805.5	0.00378
<b>Butylamine</b>	9.1	MeCN	7782.7	0.01256
<b>Haloperidol</b>	8	MeCN	7780.2	0.00215
<b>Reserpine</b>	11.6	MeCN	7771.6	0.00193
<b>Phenazine</b>	10.7	MeCN	7760.3	0.00601
<b>Clotrimazole</b>	6.8	MeCN	7776.3	0.00199
<b>Simazine</b>	2.7	MeOH	7733.1	0.00137
<b>Efavirenz</b>	10.3	MeCN	7757.7	0.00331
<b>Adenosine</b>	3.4	H <sub>2</sub> O	10001.6	0.00127
<b>Climbazole</b>	11.9	MeCN	7766.6	0.00411
<b>Melamine</b>	10.2	H <sub>2</sub> O	11982.8	0.00673
<b>Metazachlor</b>	9.2	MeCN	7758.6	0.00336
<b>Chlorothiazide</b>	14.1	MeCN	7746.4	0.00484
<b>Metformin</b>	10.7	H <sub>2</sub> O	9962.3	0.00647
<b>2-Methylbenzothiazole</b>	12.1	MeCN	7746.6	0.00823
<b>Benzothiazole</b>	13.5	MeCN	7726.9	0.01016
<b>Chlorpyrifos</b>	10.6	MeCN	7822.1	0.00304
<b>5-Methyl-1H-</b> <b>benzotriazole</b>	9.8	MeCN	7771.5	0.00744
<b>10,11-dihydro-10-</b> <b>hydroxycarbamazepine</b>	1	MeCN	-	0.00393
<b>Sudan I</b>	2.3	MeCN	9315	0.00078
<b>Ketoconazole</b>	7.3	MeCN	7807	0.00138
<b>5-Chlorobenzotriazole</b>	14.2	MeCN	7763.7	0.00936
<b>Benzotriazole-5-</b> <b>carboxylic acid</b>	9.1	MeOH:HCl 10:3	10682.3	0.00452
<b>Carbamazepine-10,11-</b> <b>epoxide</b>	1	MeOH	-	0.00396
<b>Metolachlor-ESA</b>	2.2	MeCN	7793.4	0.00063
<b>Metolachlor-OXA</b>	1.7	MeCN	7814.4	0.00061
<b>Omethoate</b>	12.4	MeCN	7772.9	0.00588
<b>Atrazine-2-hydroxy</b>	9.9	MeOH:HCl 2:3	9158.3	0.00517
<b>2-</b> <b>Hydroxybenzothiazole</b>	14.4	MeCN	7781	0.00962
<b>2-Aminobenzothiazole</b>	11	MeCN	7783.2	0.00740

Compound	$m_{\text{compound}}$ [mg]	Solvent	$m_{\text{solvent}}$ [mg]	$c$ [M]
<b>Atrazine-d<sub>5</sub></b>	3.1	MeCN	2323.4	0.00475
<b>Imazalil-d<sub>5</sub></b>	5.5	MeCN	3861.6	0.00370
<b>Caffeine-<sup>13</sup>C<sub>3</sub></b>	1	MeOH	-	0.00507

**Table S2.** The concentrations for the real samples, where SB1 & SB2 corresponds Ladviken samples and SB3 & SB4 corresponds to tap water samples.

Compound	c SB1 & SB3 [M]	c SB2 & SB4 [M]
<b>Amitrole</b>	1.03e-07	1.03e-07
<b>Aspartame</b>	4.62e-08	4.62e-08
<b>Atrazine</b>	4.07e-08	4.07e-08
<b>Benzotriazole</b>	1.24e-07	1.24e-07
<b>Butocarboxim</b>	5.44e-08	5.44e-08
<b>Caffeine</b>	6.84e-08	6.84e-08
<b>Carbamazepine</b>	3.65e-08	3.65e-08
<b>Chlormequat</b>	6.72e-08	6.72e-08
<b>Clarithromycin</b>	1.40e-08	1.40e-08
<b>Dichlorvos</b>	6.66e-08	6.66e-08
<b>Emamectin B1a</b>	1.20e-08	1.20e-08
<b>Guanylurea</b>	7.85e-08	7.85e-08
<b>Histamine</b>	9.84e-08	9.84e-08
<b>Imazalil</b>	3.31e-08	3.31e-08
<b>Irgarol</b>	2.74e-08	2.74e-08
<b>Ivermectin B1a</b>	1.29e-08	1.29e-08
<b>Methamidophos</b>	1.24e-07	1.24e-07
<b>Metolachlor</b>	4.20e-08	4.20e-08
<b>Nigericin</b>	1.32e-08	1.32e-08
<b>Octocrylene</b>	2.59e-08	2.59e-08
<b>Rifaximin</b>	1.22e-08	1.22e-08
<b>Simvastatin</b>	2.32e-08	2.32e-08
<b>Spinosyn A</b>	1.39e-08	1.39e-08
<b>Sucratose</b>	3.43e-08	3.43e-08
<b>TCMTB</b>	4.35e-08	4.34e-08
<b>Trichlorfon</b>	3.12e-08	3.12e-08
<b>Tylosin</b>	1.08e-08	1.08e-08
<b>Vancomycin</b>	5.91e-09	5.91e-09
<b>Atrazine- desethyl</b>	2.23e-07	1.11e-07
<b>Atrazine-desethyl-2-hydroxy</b>	3.18e-07	1.59e-07
<b>Atrazine-desisopropyl</b>	2.07e-07	1.04e-07
<b>Atrazine-desethyl-desisopropyl</b>	3.98e-07	1.99e-07
<b>Atrazine-desethyl-desisopropyl-2-hydroxy</b>	4.79e-07	2.40e-07
<b>Atrazine-desisopropyl-2-hydroxy</b>	3.62e-07	1.81e-07
<b>2-(Methylthio)benzothiazole</b>	2.95e-07	1.47e-07
<b>Progesterone</b>	1.36e-07	6.81e-08
<b>Butylamine</b>	1.02e-06	5.09e-07
<b>Haloperidol</b>	8.58e-08	4.29e-08
<b>Reserpine</b>	7.19e-08	3.60e-08

Compound	c SB1 & SB3 [M]	c SB2 & SB4 [M]
<b>Phenazine</b>	7.74e-08	3.87e-08
<b>Clotrimazole</b>	1.72e-07	8.59e-08
<b>Simazine</b>	2.77e-07	1.38e-07
<b>Efavirenz</b>	1.57e-06	7.85e-07
<b>Adenosine</b>	3.05e-07	1.52e-07
<b>Climbazole</b>	5.29e-08	2.65e-08
<b>Melamine</b>	9.56e-07	4.78e-07
<b>Metazachlor</b>	1.17e-07	5.83e-08
<b>Chlorothiazide</b>	2.30e-06	1.15e-06
<b>Metformin</b>	5.97e-07	2.98e-07
<b>2-Methylbenzothiazole</b>	7.31e-07	3.65e-07
<b>Benzothiazole</b>	4.90e-06	2.45e-06
<b>Chlorpyrifos</b>	2.94e-06	1.47e-06
<b>5-Methyl-1H-benzotriazole</b>	2.78e-07	1.39e-07
<b>10,11-dihydro-10-hydroxycarbamazepine</b>	1.32e-07	6.58e-08
<b>Sudan I</b>	1.04e-07	5.18e-08
<b>Ketoconazole</b>	6.50e-07	3.25e-07
<b>5-Chlorobenzotriazole</b>	7.83e-07	3.92e-07
<b>Benzotriazole-5-carboxylic acid</b>	2.17e-06	1.08e-06
<b>Carbamazepine-10,11-epoxide</b>	7.09e-08	3.54e-08
<b>Metolachlor-ESA</b>	1.56e-06	7.79e-07
<b>Metolachlor-OXA</b>	1.49e-06	7.47e-07
<b>Omethoate</b>	2.57e-07	1.29e-07
<b>Atrazine-2-hydroxy</b>	9.98e-08	4.99e-08
<b>2-Hydroxybenzothiazole</b>	3.96e-07	1.98e-07
<b>2-Aminobenzothiazole</b>	2.66e-07	1.33e-07
<b>Atrazine-d<sub>5</sub></b>	4.15e-08	4.15e-08
<b>Caffeine-<sup>13</sup>C<sub>3</sub></b>	3.36e-08	3.36e-08
<b>Imazalil-d<sub>5</sub></b>	3.09e-08	3.09e-08



**Figure S1.** The stability plots with data from week 1 included. As seen, something happened with the instrumental stability in week 1. This is clearest in the left plot.

**Table S3.** Comparison of the predicted concentrations for all approaches and the real concentration for all suspect compounds in all four samples.

Compound	Sample	c <sub>real</sub>	Closest Eluting Std c <sub>pred</sub>	Ionisation Efficiency c <sub>pred</sub>	Parent Compound c <sub>pred</sub>	Similar Compound c <sub>pred</sub>
<b>10,11-Dihydro-10-hydroxycarbamazepine</b>	SB1	1.32e-07	3.22e-08	7.04e-07	NA	NA
<b>2-(Methylthio)benzothiazole</b>	SB1	2.95e-07	8.25e-07	1.98e-07	8.25e-07	8.25e-07
<b>2-Aminobenzothiazole</b>	SB1	2.67e-07	1.97e-06	2.28e-06	6.13e-06	6.13e-06
<b>2-Hydroxybenzothiazole</b>	SB1	3.96e-07	8.54e-07	3.39e-08	4.79e-08	1.05e-08
<b>2-Methylbenzothiazole</b>	SB1	7.31e-07	4.18e-08	4.07e-07	NA	NA
<b>5-Chlorobenzotriazole</b>	SB1	7.83e-07	9.55e-08	1.07e-06	4.05e-07	4.05e-07
<b>5-Methyl-1H-benzotriazole</b>	SB1	2.78e-07	5.48e-08	1.73e-06	1.03e-06	1.03e-06
<b>Adenosine</b>	SB1	3.05e-07	1.36e-06	1.57e-06	NA	2.78e-08
<b>Atrazine-2-hydroxy</b>	SB1	9.98e-08	1.45e-06	1.06e-06	1.42e-07	1.42e-07
<b>Atrazine-desethyl</b>	SB1	2.23e-07	4.81e-08	7.69e-07	1.29e-07	1.29e-07
<b>Atrazine-desethyl-2-hydroxy</b>	SB1	3.18e-07	2.60e-06	2.01e-06	1.43e-07	1.43e-07
<b>Atrazine-desethyl-desisopropyl</b>	SB1	3.98e-07	1.76e-07	5.06e-07	2.99e-08	9.72e-07
<b>Atrazine-desethyl-desisopropyl-2-hydroxy</b>	SB1	4.79e-07	5.64e-08	7.89e-08	3.10e-09	1.01e-07
<b>Atrazine-desisopropyl</b>	SB1	2.07e-07	1.48e-04	5.48e-07	8.14e-08	8.14e-08
<b>Atrazine-desisopropyl-2-hydroxy</b>	SB1	3.62e-07	4.36e-06	2.73e-06	2.39e-07	2.39e-07
<b>Benzothiazole</b>	SB1	4.90e-06	2.43e-05	2.49e-06	4.92e-06	1.08e-06
<b>Benzotriazole-5-carboxylic acid</b>	SB1	2.17e-06	1.78e-07	2.32e-06	NA	NA
<b>Carbamazepine-10,11-epoxide</b>	SB1	7.09e-08	1.33e-04	1.38e-06	NA	NA
<b>Chlorothiazide</b>	SB1	2.30e-06	3.09e-08	3.74e-07	NA	4.29e-09
<b>Chlorpyrifos</b>	SB1	2.94e-06	2.40e-08	2.35e-08	NA	5.48e-09

Compound	Sample	c <sub>real</sub>	Closest Eluting Std c <sub>pred</sub>	Ionisation Efficiency c <sub>pred</sub>	Parent Compound c <sub>pred</sub>	Similar Compound c <sub>pred</sub>
<b>Climbazole</b>	SB1	5.29e-08	1.64e-07	5.25e-07	NA	8.90e-08
<b>Clotrimazole</b>	SB1	1.72e-07	5.13e-07	1.20e-06	NA	3.96e-07
<b>Efavirenz</b>	SB1	1.57e-06	7.48e-09	5.19e-08	NA	1.90e-09
<b>Haloperidol</b>	SB1	8.58e-08	8.79e-08	1.21e-07	NA	8.79e-08
<b>Ketoconazole</b>	SB1	6.50e-07	2.00e-06	1.02e-06	NA	4.52e-07
<b>Metazachlor</b>	SB1	1.17e-07	1.88e-07	1.03e-06	NA	3.94e-07
<b>Metformin</b>	SB1	5.97e-07	2.25e-05	4.74e-06	NA	1.21e-05
<b>Omethoate</b>	SB1	2.57e-07	4.46e-06	5.75e-06	NA	6.86e-05
<b>Phenazine</b>	SB1	7.74e-08	4.94e-08	1.81e-07	NA	2.63e-08
<b>Progesterone</b>	SB1	1.36e-07	2.84e-08	1.32e-07	NA	2.84e-08
<b>Reserpine</b>	SB1	7.19e-08	4.04e-08	4.86e-08	NA	8.60e-08
<b>Simazine</b>	SB1	2.77e-07	4.05e-05	6.68e-07	NA	2.57e-07
<b>Sudan I</b>	SB1	1.04e-07	2.19e-07	2.57e-08	NA	1.81e-08
<b>10,11-Dihydro-10-hydroxycarbamazepine</b>	SB2	6.58e-08	1.65e-08	3.38e-07	NA	NA
<b>2-(Methylthio)benzothiazole</b>	SB2	1.47e-07	3.90e-07	9.54e-08	3.90e-07	3.90e-07
<b>2-Aminobenzothiazole</b>	SB2	1.33e-07	9.05e-07	1.05e-06	2.77e-06	2.77e-06
<b>2-Hydroxybenzothiazole</b>	SB2	1.98e-07	4.04e-07	1.30e-08	1.82e-08	4.24e-09
<b>2-Methylbenzothiazole</b>	SB2	3.65e-07	2.11e-08	1.92e-07	NA	NA
<b>5-Chlorobenzotriazole</b>	SB2	3.92e-07	5.62e-08	5.69e-07	2.25e-07	2.25e-07
<b>5-Methyl-1H-benzotriazole</b>	SB2	1.39e-07	2.81e-08	8.32e-07	5.15e-07	5.15e-07
<b>Adenosine</b>	SB2	1.52e-07	7.85e-07	8.39e-07	NA	1.59e-08
<b>Atrazine-2-hydroxy</b>	SB2	4.99e-08	7.35e-07	5.40e-07	7.64e-08	7.64e-08
<b>Atrazine-desethyl</b>	SB2	1.11e-07	1.90e-08	NA	5.07e-08	5.07e-08
<b>Atrazine-desethyl-2-hydroxy</b>	SB2	1.59e-07	1.39e-06	9.94e-07	7.53e-08	7.53e-08

Compound	Sample	c <sub>real</sub>	Closest Eluting Std c <sub>pred</sub>	Ionisation Efficiency c <sub>pred</sub>	Parent Compound c <sub>pred</sub>	Similar Compound c <sub>pred</sub>
<b>Atrazine-desethyl-</b> <b>desisopropyl</b>	SB2	1.99e-07	9.14e-08	2.36e-07	1.49e-08	4.35e-07
<b>Atrazine-desethyl-</b> <b>desisopropyl-2-hydroxy</b>	SB2	2.40e-07	2.74e-08	3.53e-08	1.48e-09	4.34e-08
<b>Atrazine-desisopropyl</b>	SB2	1.04e-07	3.45e-05	2.60e-07	4.09e-08	4.09e-08
<b>Atrazine-desisopropyl-2-</b> <b>hydroxy</b>	SB2	1.81e-07	2.38e-06	1.39e-06	1.29e-07	1.29e-07
<b>Benzothiazole</b>	SB2	2.45e-06	1.07e-05	1.16e-06	2.27e-06	5.28e-07
<b>Benzotriazole-5-carboxylic</b> <b>acid</b>	SB2	1.08e-06	9.31e-08	1.12e-06	NA	NA
<b>Carbamazepine-10,11-</b> <b>epoxide</b>	SB2	3.54e-08	8.09e-05	6.89e-07	NA	NA
<b>Chlorothiazide</b>	SB2	1.15e-06	1.47e-08	1.64e-07	NA	1.92e-09
<b>Chlorpyrifos</b>	SB2	1.47e-06	3.29e-09	1.02e-08	NA	2.48e-09
<b>Climbazole</b>	SB2	2.65e-08	8.11e-08	2.56e-07	NA	4.55e-08
<b>Clotrimazole</b>	SB2	8.59e-08	2.53e-07	5.84e-07	NA	1.90e-07
<b>Efavirenz</b>	SB2	7.85e-07	3.49e-09	2.32e-08	NA	9.05e-10
<b>Haloperidol</b>	SB2	4.29e-08	4.52e-08	6.00e-08	NA	4.52e-08
<b>Ketoconazole</b>	SB2	3.25e-07	1.01e-06	4.80e-07	NA	2.20e-07
<b>Metazachlor</b>	SB2	5.83e-08	9.55e-08	4.63e-07	NA	1.83e-07
<b>Metformin</b>	SB2	2.98e-07	1.27e-05	2.81e-06	NA	7.52e-06
<b>Omethoate</b>	SB2	1.29e-07	2.22e-06	2.85e-06	NA	3.22e-05
<b>Phenazine</b>	SB2	3.87e-08	2.56e-08	8.89e-08	NA	1.28e-08
<b>Progesterone</b>	SB2	6.81e-08	1.46e-08	6.54e-08	NA	1.46e-08
<b>Reserpine</b>	SB2	3.60e-08	1.78e-08	2.14e-08	NA	3.86e-08
<b>Simazine</b>	SB2	1.38e-07	1.85e-05	3.28e-07	NA	1.32e-07
<b>Sudan I</b>	SB2	5.18e-08	3.05e-08	1.16e-08	NA	8.01e-09

Compound	Sample	c <sub>real</sub>	Closest Eluting Std	Ionisation Efficiency	Parent Compound	Similar Compound
		c <sub>pred</sub>	c <sub>pred</sub>	c <sub>pred</sub>	c <sub>pred</sub>	c <sub>pred</sub>
<b>10,11-Dihydro-10-hydroxycarbamazepine</b>	SB3	1.32e-07	3.40e-08	7.01e-07	NA	NA
<b>2-(Methylthio)benzothiazole</b>	SB3	2.95e-07	8.03e-07	1.81e-07	8.74e-07	8.74e-07
<b>2-Aminobenzothiazole</b>	SB3	2.67e-07	1.97e-06	2.18e-06	6.67e-06	6.67e-06
<b>2-Hydroxybenzothiazole</b>	SB3	3.96e-07	7.32e-07	2.84e-08	4.47e-08	8.43e-09
<b>2-Methylbenzothiazole</b>	SB3	7.31e-07	4.23e-08	3.85e-07	NA	NA
<b>5-Chlorobenzotriazole</b>	SB3	7.83e-07	1.02e-07	1.11e-06	4.05e-07	4.05e-07
<b>5-Methyl-1H-benzotriazole</b>	SB3	2.78e-07	5.42e-08	1.64e-06	9.51e-07	9.51e-07
<b>Adenosine</b>	SB3	3.05e-07	1.35e-06	1.53e-06	NA	2.59e-08
<b>Atrazine-2-hydroxy</b>	SB3	9.98e-08	1.49e-06	1.02e-06	1.25e-07	1.25e-07
<b>Atrazine-desethyl</b>	SB3	2.23e-07	3.99e-08	NA	9.20e-08	9.20e-08
<b>Atrazine-desethyl-2-hydroxy</b>	SB3	3.18e-07	2.36e-06	1.76e-06	1.12e-07	1.12e-07
<b>Atrazine-desethyl-desisopropyl</b>	SB3	3.98e-07	1.73e-07	4.83e-07	2.54e-08	9.48e-07
<b>Atrazine-desethyl-desisopropyl-2-hydroxy</b>	SB3	4.79e-07	4.81e-08	6.62e-08	2.28e-09	8.51e-08
<b>Atrazine-desisopropyl</b>	SB3	2.07e-07	1.54e-04	5.33e-07	7.27e-08	7.27e-08
<b>Atrazine-desisopropyl-2-hydroxy</b>	SB3	3.62e-07	4.33e-06	2.59e-06	2.05e-07	2.05e-07
<b>Benzothiazole</b>	SB3	4.90e-06	2.61e-05	2.59e-06	5.76e-06	1.09e-06
<b>Benzotriazole-5-carboxylic acid</b>	SB3	2.17e-06	1.83e-07	2.39e-06	NA	NA
<b>Carbamazepine-10,11-epoxide</b>	SB3	7.09e-08	1.36e-04	1.32e-06	NA	NA
<b>Chlorothiazide</b>	SB3	2.30e-06	3.52e-08	4.27e-07	NA	4.56e-09
<b>Chlorpyrifos</b>	SB3	2.94e-06	2.11e-08	1.92e-08	NA	4.37e-09
<b>Climbazole</b>	SB3	5.29e-08	1.75e-07	5.04e-07	NA	8.83e-08

Compound	Sample	Closest Eluting Std		Ionisation Efficiency	Parent Compound	Similar Compound
		c <sub>real</sub>	c <sub>pred</sub>			
<b>Clotrimazole</b>	SB3	1.72e-07	5.65e-07	1.17e-06	NA	4.07e-07
<b>Efavirenz</b>	SB3	1.57e-06	7.62e-09	5.04e-08	NA	1.81e-09
<b>Haloperidol</b>	SB3	8.58e-08	9.59e-08	1.13e-07	NA	8.95e-08
<b>Ketoconazole</b>	SB3	6.50e-07	2.03e-06	9.08e-07	NA	4.29e-07
<b>Metazachlor</b>	SB3	1.17e-07	2.01e-07	1.02e-06	NA	3.83e-07
<b>Metformin</b>	SB3	5.97e-07	2.09e-05	4.12e-06	NA	1.12e-05
<b>Omethoate</b>	SB3	2.57e-07	4.36e-06	5.39e-06	NA	6.42e-05
<b>Phenazine</b>	SB3	7.74e-08	5.00e-08	1.68e-07	NA	2.49e-08
<b>Progesterone</b>	SB3	1.36e-07	3.06e-08	1.34e-07	NA	2.78e-08
<b>Reserpine</b>	SB3	7.19e-08	4.17e-08	4.37e-08	NA	8.07e-08
<b>Simazine</b>	SB3	2.77e-07	4.17e-05	6.23e-07	NA	2.28e-07
<b>Sudan I</b>	SB3	1.04e-07	2.11e-07	2.15e-08	NA	1.63e-08
<b>10,11-Dihydro-10-hydroxycarbamazepine</b>	SB4	6.58e-08	1.80e-08	5.17e-07	NA	NA
<b>2-(Methylthio)benzothiazole</b>	SB4	1.47e-07	3.85e-07	1.21e-07	4.18e-07	4.18e-07
<b>2-Aminobenzothiazole</b>	SB4	1.33e-07	9.79e-07	1.52e-06	3.31e-06	3.31e-06
<b>2-Hydroxybenzothiazole</b>	SB4	1.98e-07	3.58e-07	1.97e-08	2.18e-08	4.29e-09
<b>2-Methylbenzothiazole</b>	SB4	3.65e-07	2.16e-08	2.74e-07	NA	NA
<b>5-Chlorobenzotriazole</b>	SB4	3.92e-07	4.26e-08	6.54e-07	1.76e-07	1.76e-07
<b>5-Methyl-1H-benzotriazole</b>	SB4	1.39e-07	2.50e-08	1.06e-06	4.56e-07	4.56e-07
<b>Adenosine</b>	SB4	1.52e-07	7.02e-07	1.14e-06	NA	1.61e-08
<b>Atrazine-2-hydroxy</b>	SB4	4.99e-08	7.61e-07	7.31e-07	6.48e-08	6.48e-08
<b>Atrazine-desethyl</b>	SB4	1.11e-07	1.93e-08	NA	4.52e-08	4.52e-08
<b>Atrazine-desethyl-2-hydroxy</b>	SB4	1.59e-07	1.24e-06	1.29e-06	5.92e-08	5.92e-08
<b>Atrazine-desethyl-desisopropyl</b>	SB4	1.99e-07	7.90e-08	3.11e-07	1.17e-08	3.88e-07

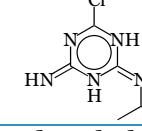
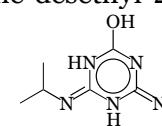
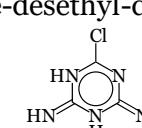
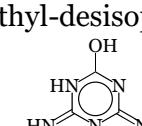
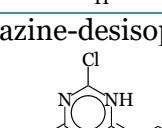
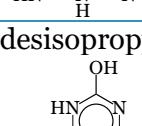
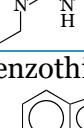
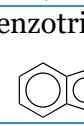
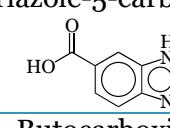
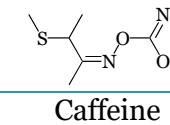
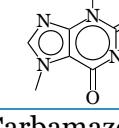
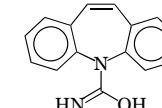
Compound	Sample	c <sub>real</sub>	Closest Eluting Std	Ionisation Efficiency	Parent Compound	Similar Compound
			c <sub>pred</sub>	c <sub>pred</sub>	c <sub>pred</sub>	c <sub>pred</sub>
<b>Atrazine-desethyl-desisopropyl-2-hydroxy</b>	SB4	2.40e-07	2.46e-08	4.79e-08	1.18e-09	3.89e-08
<b>Atrazine-desisopropyl</b>	SB4	1.04e-07	7.81e-05	3.77e-07	3.74e-08	3.74e-08
<b>Atrazine-desisopropyl-2-hydroxy</b>	SB4	1.81e-07	2.36e-06	1.98e-06	1.13e-07	1.13e-07
<b>Benzothiazole</b>	SB4	2.45e-06	1.30e-05	1.82e-06	2.86e-06	5.63e-07
<b>Benzotriazole-5-carboxylic acid</b>	SB4	1.08e-06	8.57e-08	1.60e-06	NA	NA
<b>Carbamazepine-10,11-epoxide</b>	SB4	3.54e-08	6.68e-05	8.99e-07	NA	NA
<b>Chlorothiazide</b>	SB4	1.15e-06	1.50e-08	2.60e-07	NA	1.89e-09
<b>Chlorpyrifos</b>	SB4	1.47e-06	7.43e-09	9.43e-09	NA	1.70e-09
<b>Climbazole</b>	SB4	2.65e-08	8.16e-08	3.23e-07	NA	4.93e-08
<b>Clotrimazole</b>	SB4	8.59e-08	1.83e-07	5.19e-07	NA	1.28e-07
<b>Efavirenz</b>	SB4	7.85e-07	2.75e-09	2.55e-08	NA	7.82e-10
<b>Haloperidol</b>	SB4	4.29e-08	3.96e-08	6.27e-08	NA	4.42e-08
<b>Ketoconazole</b>	SB4	3.25e-07	6.77e-07	4.10e-07	NA	1.71e-07
<b>Metazachlor</b>	SB4	5.83e-08	9.64e-08	6.74e-07	NA	2.03e-07
<b>Metformin</b>	SB4	2.98e-07	1.33e-05	3.66e-06	NA	7.61e-06
<b>Omethoate</b>	SB4	1.29e-07	2.25e-06	3.91e-06	NA	3.41e-05
<b>Phenazine</b>	SB4	3.87e-08	2.55e-08	1.18e-07	NA	1.23e-08
<b>Progesterone</b>	SB4	6.81e-08	1.45e-08	8.85e-08	NA	1.46e-08
<b>Reserpine</b>	SB4	3.60e-08	1.42e-08	2.01e-08	NA	3.03e-08
<b>Simazine</b>	SB4	1.38e-07	2.07e-05	4.26e-07	NA	1.15e-07
<b>Sudan I</b>	SB4	5.18e-08	6.50e-08	8.97e-09	NA	4.87e-09

**Table S4.** All suspect compound displayed together with their respective closest eluting-, parent-, and most similar compound.

Suspect compound	Closest eluting compound	Parent compound	Most similar compound
<b>10,11-Dihydro-10-hydroxycarbamazepine</b>	Imazalil	NA	NA
<b>2-(Methylthio)benzothiazole</b>	TCMTB	TCMTB	TCMTB
<b>2-Aminobenzothiazole</b>	Vancomycin [M] <sup>2+</sup>	TCMTB	TCMTB
<b>2-Hydroxybenzothiazole</b>	Butocarboxim	TCMTB	Benzotriazole
<b>2-Methylbenzothiazole</b>	Spinosyn A	NA	NA
<b>5-Chlorobenzotriazole</b>	Tylosin	Benzotriazole	Benzotriazole
<b>5-Methyl-1H-benzotriazole</b>	Imazalil	Benzotriazole	Benzotriazole
<b>Adenosine</b>	Amitrole	NA	Imazalil
<b>Atrazine-2-hydroxy</b>	Vancomycin [M] <sup>2+</sup>	Atrazine	Atrazine
<b>Atrazine-desethyl</b>	Imazalil	Atrazine	Atrazine
<b>Atrazine-desethyl-2-hydroxy</b>	Amitrole	Atrazine	Atrazine
<b>Atrazine-desethyl-desisopropyl</b>	Methamidophos	Atrazine	Guanlylurea
<b>Atrazine-desethyl-desisopropyl-2-hydroxy</b>	Amitrole	Atrazine	Guanlylurea
<b>Atrazine-desisopropyl</b>	Sucralose [M+Na] <sup>+</sup>	Atrazine	Atrazine
<b>Atrazine-desisopropyl-2-hydroxy</b>	Amitrole	Atrazine	Atrazine
<b>Benzothiazole</b>	Dichlorvos	TCMTB	Benzotriazole
<b>Benzotriazole-5-carboxylic acid</b>	Aspartame	NA	NA
<b>Carbamazepine-10,11-epoxide</b>	Butocarboxim	NA	NA
<b>Chlorothiazide</b>	Aspartame	NA	Carbamazepine
<b>Chlorpyrifos</b>	Simvastatin	NA	Metolachlor
<b>Climbazole</b>	Clarithromycin	NA	Imazalil
<b>Clotrimazole</b>	Clarithromycin	NA	Carbamazepine
<b>Efavirenz</b>	Metolachlor	NA	Imazalil
<b>Haloperidol</b>	Imazalil	NA	Imazalil
<b>Ketoconazole</b>	Tylosin	NA	Imazalil
<b>Metazachlor</b>	Emamectin B1a	NA	Metolachlor
<b>Metformin</b>	Guanlylurea	NA	Caffeine
<b>Omethoate</b>	Vancomycin [M] <sup>2+</sup>	NA	Dichlorvos
<b>Phenazine</b>	Atrazine	NA	Carbamazepine
<b>Progesterone</b>	Metolachlor	NA	Metolachlor
<b>Reserpine</b>	Clarithromycin	NA	Metolachlor
<b>Simazine</b>	Dichlorvos	NA	Atrazine
<b>Sudan I</b>	Simvastatin	NA	Carbamazepine

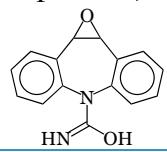
**Table S5.** The structures of all the compounds that were visible in the analysis.

Structure
10,11-Dihydro-10-hydroxycarbamazepine 
2-(Methylthio)benzothiazole 
2-Aminobenzothiazole 
2-Hydroxybenzothiazole 
2-Methylbenzothiazole 
5-Chlorobenzotriazole 
5-Methyl-1H-benzotriazole 
Adenosine 
Amitrole 
Aspartame 
Atrazine 
Atrazine-2-hydroxy 

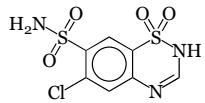
Structure
Atrazine-desethyl 
Atrazine-desethyl-2-hydroxy 
Atrazine-desethyl-desisopropyl 
Atrazine-desethyl-desisopropyl-2-hydroxy 
Atrazine-desisopropyl 
Atrazine-desisopropyl-2-hydroxy 
Benzothiazole 
Benzotriazole 
Benzotriazole-5-carboxylic acid 
Butocarboxim 
Caffeine 
Carbamazepine 

## Structure

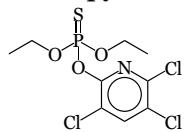
Carbamazepine-10,11-epoxide



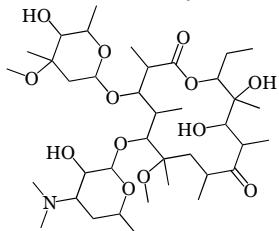
Chlorothiazide



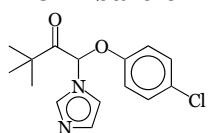
Chlorpyrifos



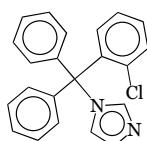
Clarithromycin



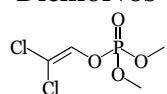
Climbazole



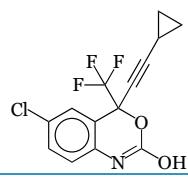
Clotrimazole



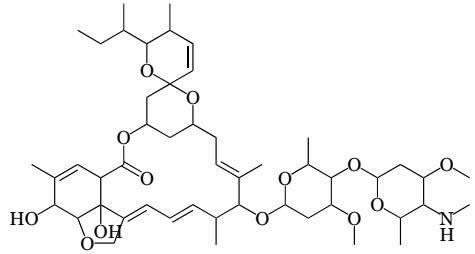
Dichlorvos

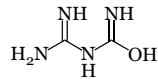
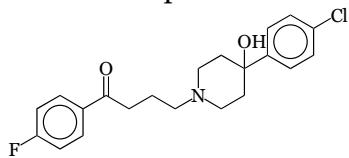
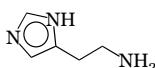
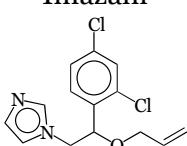
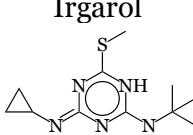
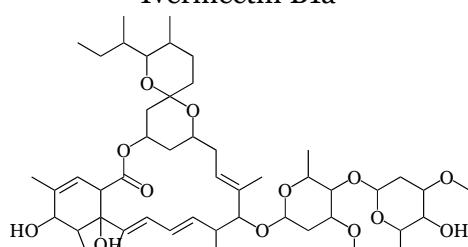
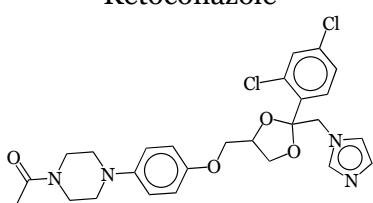
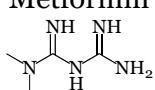


Efavirenz



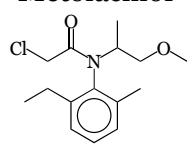
Emamectin B1a



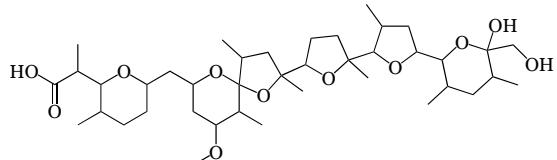
Structure
Guanlylurea 
Haloperidol 
Histamine 
Imazalil 
Irgarol 
Ivermectin B1a 
Ketoconazole 
Metazachlor 
Metformin 
Methamidophos 

## Structure

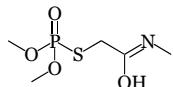
Metolachlor



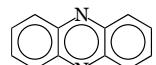
Nigericin



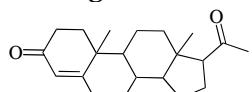
Omethoate



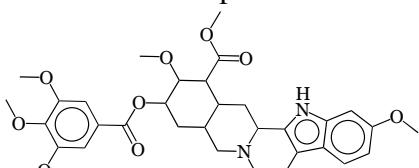
Phenazine



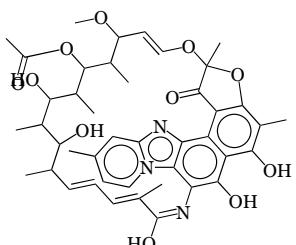
Progesterone



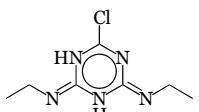
Reserpine



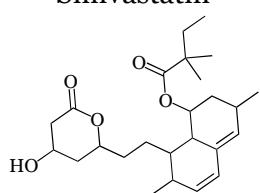
Rifaximin



Simazine

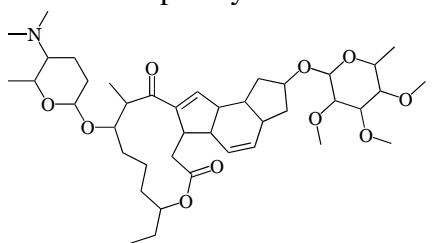


Simvastatin

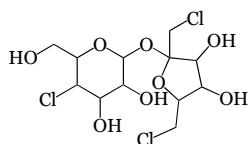


## Structure

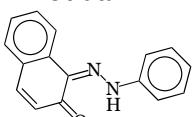
Spinosyn A



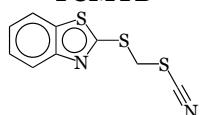
Sucratose



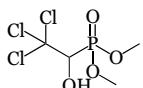
Sudan I



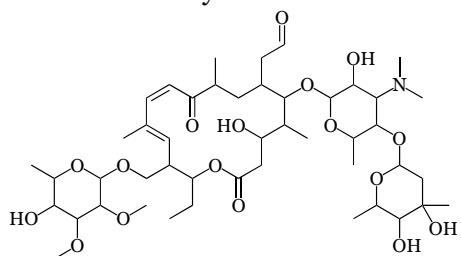
TCMTB



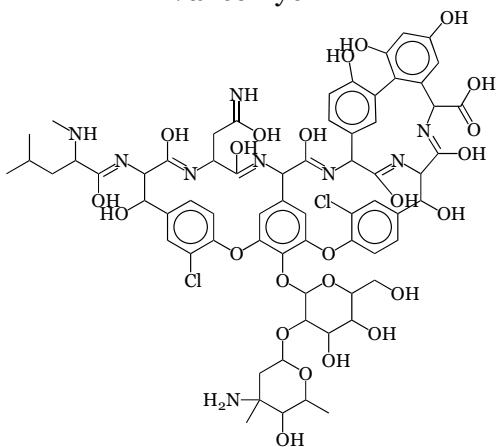
Trichlorfon



Tylosin



Vancomycin



**Code 1.** The code used for the semi-quantification methods.

```
library(caret)
library(enviPat)
library(plotly)
library(RRF)
library(rcdk)
library(tidyverse)
library(webchem)

fn_isotope_distribution <- function(smiles) {
  molecule <- parse.smiles(smiles)[[1]]
  formula <- get.mol2formula(molecule, charge=0)
  formula <- formula@string

  data(isotopes)
  pattern<-isopattern(isotopes,
                        formula,
                        threshold=0.1,
                        plotit=FALSE,
                        charge=FALSE,
                        algo=1)
  isotopes <- as.data.frame(pattern[[1]])
  isotope_dist <- as.numeric(sum(isotopes$abundance))
  return(isotope_dist)
}

fn_viscosity <- function(organic, organic_modifier){
  viscosity <- case_when(
    organic_modifier == "MeCN" ~ (-
0.000103849885417527)*organic^2+0.00435719229180079*organic+0.884232
851261593,
    organic_modifier == "MeOH" ~ (-
0.00035908)*organic^2+0.031972067*organic+0.90273943)
  return(viscosity)
}

fn_surface_tension <- function(organic, organic_modifier){
  surface_tension <- case_when(
    organic_modifier == "MeCN" ~ 71.76-
2.906*71.76*(organic/100)+(7.138*27.86+2.906*71.76-
71.76)*(organic/100)^2+(27.86-7.138*27.86)*(organic/100)^3,
    organic_modifier == "MeOH" ~ 71.76-
2.245*71.76*(organic/100)+(5.625*22.12+2.245*71.76-
71.76)*(organic/100)^2+(22.12-5.625*22.12)*(organic/100)^3)
  return(surface_tension)
}

fn_polarity_index <- function(organic, organic_modifier){
  polarity_index <- case_when(
    organic_modifier == "MeCN" ~ (organic/100)*5.1+((100-
organic)/100)*10.2,
    organic_modifier == "MeOH" ~ (organic/100)*5.1+((100-
organic)/100)*10.2)
  return(polarity_index)
}
```

```

fn_organic_percentage <- function(eluent_parameters, ret_time) {
  ApproxFun <- approxfun(x = eluent_parameters$time, y =
eluent_parameters$B)
  organic <- ApproxFun(ret_time)
  return(organic)
}

assigning_closest_cal_comp_RF <- function(xRT, cal_compounds) {
  RF_cal <- cal_compounds %>% slice(which.min(abs(xRT - ret_time)))
%>%select(RF)
  print(unlist(RF_cal))
}

assigning_closest_cal_comp <- function(xRT, cal_compounds) {
  Comp_cal <- cal_compounds %>% slice(which.min(abs(xRT -
ret_time))) %>%select(Compound)
  print(unlist(Comp_cal))
}

assigning_closest_cal_comp_RT <- function(xRT, cal_compounds) {
  Comp_cal <- cal_compounds %>% slice(which.min(abs(xRT -
ret_time))) %>%select(ret_time)
  print(unlist(Comp_cal))
}

SB1 <- read_delim('Quantem_SB1_w_compname.csv',
                   delim = ';',
                   col_names = TRUE)

SB1 <- SB1 %>%
  group_by(SMILES) %>%
  mutate(IC = fn_isotope_distribution(SMILES)) %>%
  ungroup()

SB1 <- SB1 %>%
  mutate(signal = signal*IC,
        RF = signal / concentration)

SA <- SB1 %>%
  filter(`SB/SA comp.` == "SA")

similarity_matrix <- read_delim('similarity_matrix.csv',
                                 delim = ";",
                                 col_names = TRUE)

similarity_matrix <- similarity_matrix %>%
  rename(SB_compound = Compound) %>%
  rename(Compound = Similar) %>%
  group_by(SB_compound) %>%
  slice(which.max(similarity_per_cent)) %>%
  ungroup() %>%
  left_join(SA) %>%

```

```

  select(SB_compound, Compound, RF, ret_time, similarity_per_cent,
Maximum_expected_error) %>%
  rename(Similar = Compound) %>%
  rename(Compound = SB_compound,
     RF_similar_compound = RF,
     ret_time_similar_compound = ret_time)

SB1 <- SB1 %>%
  left_join(similarity_matrix) %>%
  mutate(c_pred_similar_compound = signal / RF_similar_compound)

cal_Comp_closeRT <- c()
for(i in 1:69){
  cal_Comp_closeRT <- c(cal_Comp_closeRT,
assigning_closest_cal_comp(SB1[i,]$ret_time, SA))
}

cal_RF_closeRT <- c()
for(i in 1:69){
  cal_RF_closeRT <- c(cal_RF_closeRT,
assigning_closest_cal_comp_RF(SB1[i,]$ret_time, SA))
}

cal_RT_closeRT <- c()
for(i in 1:69){
  cal_RT_closeRT <- c(cal_RT_closeRT,
assigning_closest_cal_comp_RT(SB1[i,]$ret_time, SA))
}

SB1 <- data.frame(SB1, cal_Comp_closeRT, cal_RF_closeRT,
cal_RT_closeRT)

SB1 <- SB1 %>%
  mutate(c_pred_close_ret_time = signal / cal_RF_closeRT)

regressor_pos <- read_rds("ESIpos_model_191116.rds")
descs_pos <- read_rds("ESIpos_model_descs_191116.rds")

Padel_data <- read_delim('descs200629.csv',
                           delim = ",",
                           col_names = TRUE,
                           trim_ws = TRUE)

eluent_parameters <- read_delim('eluent.csv',
                                 delim = ";",
                                 col_names = TRUE)

organic_modifier <- "MeCN"
pH <- 2.7
NH4 <- 0

SB1 <- SB1 %>%
  mutate(
    organic_modifier = organic_modifier,
    organic = fn_organic_percentage(eluent_parameters, ret_time),
    pH.aq. = pH,

```

```

NH4 = NH4,
viscosity = fn_viscosity(organic, organic_modifier),
surface_tension = fn_surface_tension(organic, organic_modifier),
polarity_index = fn_polarity_index(organic, organic_modifier))

prediction_set_model_pos <- SB1 %>%
  left_join(Padel_data) %>%
  select(Compound, SMILES, descs_pos) %>%
  na.omit() %>%
  mutate(logIE_pred = 0)
prediction <- predict(regressor_pos, newdata =
  prediction_set_model_pos, predict.all = TRUE)
prediction <- prediction$aggregate
prediction_set_model_pos <- prediction_set_model_pos %>%
  mutate(logIE_pred = prediction) %>%
  select(Compound, SMILES, logIE_pred)

SB1 <- SB1 %>%
  left_join(prediction_set_model_pos)

lin_fit_logRF <- lm(log(RF, 10) ~ logIE_pred,
  data = SB1 %>% filter(SB.SA.comp. == "SA") %>% filter(Compound != "Butocarboxim [M+NH4]+" &
  Compound != "Clarithromycin [M+Na]+" &
  Compound != "Ivermectin [M+NH4]+" &
  Compound != "Nigericin [M+NH4]+" &
  Compound != "Simvastatin [M+Na]+" &
  Compound != "Simvastatin [M+NH4]+" &
  Compound != "Sucralose [M+Na]+" &
  Compound != "Sucralose [M+NH4]+"))
&

SB1 <- SB1 %>%
  mutate(logRF_pred = lin_fit_logRF$coefficients[2]*logIE_pred +
  lin_fit_logRF$coefficients[1],
  c_pred_IE = signal/(10^logRF_pred))

similarity_matrix <- read_delim('similarity_matrix.csv',
  delim = ";",
  col_names = TRUE)

similarity_matrix <- similarity_matrix %>%
  rename(SB_compound = Compound) %>%
  rename(Compound = Parent_compound) %>%
  group_by(SB_compound) %>%
  slice(which.max(similarity_per_cent)) %>%
  ungroup() %>%
  left_join(SA) %>%
  select(SB_compound, Compound, RF, ret_time, similarity_per_cent,
Maximum_expected_error) %>%
  rename(Parent_compound = Compound) %>%
  rename(Compound = SB_compound,
  RF_parent_compound = RF,
  ret_time_parent_compound = ret_time)

SB1 <- SB1 %>%

```

```

left_join(similarity_matrix) %>%
  mutate(c_pred_parent_compound = signal / RF_parent_compound)

SB1 <- SB1%>%
  mutate(similar_c_comparison = case_when(concentration <
c_pred_similar_compound ~ c_pred_similar_compound / concentration,
                                             concentration >
c_pred_similar_compound ~ concentration / c_pred_similar_compound),
         closeRT_c_comparison = case_when(concentration <
c_pred_close_ret_time ~ c_pred_close_ret_time / concentration,
                                             concentration >
c_pred_close_ret_time ~ concentration / c_pred_close_ret_time),
         IE_c_comparison = case_when(concentration < c_pred_IE ~
c_pred_IE / concentration,
                                         concentration > c_pred_IE ~
concentration / c_pred_IE),
         parent_c_comparison = case_when(concentration <
c_pred_parent_compound ~ c_pred_parent_compound / concentration,
                                             concentration >
c_pred_parent_compound ~ concentration / c_pred_parent_compound))

all_c_comparison_SB1 <- SB1%>%
  select(Compound, concentration, SB.SA.comp.,
c_pred_similar_compound,
         c_pred_close_ret_time, c_pred_IE, c_pred_parent_compound,
         similar_c_comparison, closeRT_c_comparison,
IE_c_comparison,
         parent_c_comparison, ret_time, RF, similarity_per_cent)

all_c_comparison_SB1_error <- gather(data = all_c_comparison_SB1,
                                      key = "SQ_approach",
                                      value = error,
similar_c_comparison,
                                      closeRT_c_comparison,
IE_c_comparison,
                                      parent_c_comparison) %>%
  select(Compound, concentration, SB.SA.comp., ret_time, RF, error,
SQ_approach, similarity_per_cent)

all_c_comparison_SB1_pred_c <- gather(data = all_c_comparison_SB1,
                                      key = "SQ_approach",
                                      value = pred_conc,
c_pred_similar_compound,
                                      c_pred_close_ret_time,
c_pred_IE, c_pred_parent_compound, similarity_per_cent) %>%
  mutate(SQ_approach = case_when(
    SQ_approach == "c_pred_similar_compound" ~
"similar_c_comparison",
    SQ_approach == "c_pred_close_ret_time" ~ "closeRT_c_comparison",
    SQ_approach == "c_pred_IE" ~ "IE_c_comparison",
    SQ_approach == "c_pred_parent_compound" ~ "parent_c_comparison"
  )) %>%
  select(Compound, concentration, SB.SA.comp., ret_time, RF,
pred_conc, SQ_approach)

```

```

all_c_comparison_SB1 <- all_c_comparison_SB1_error %>%
  left_join(all_c_comparison_SB1_pred_c)

all_c_comparison_SB1 <- all_c_comparison_SB1 %>%
  filter(SB.SA.comp. == "SB") %>%
  filter(error != "NA")

all_c_comparison_SB1 <- all_c_comparison_SB1 %>%
  mutate(sample = "SB1")

all_c_comparison_SB1 <- all_c_comparison_SB1 %>%
  select(Compound, concentration, ret_time, RF, similarity_per_cent,
SQ_approach, error,
        pred_conc, sample)

error_stat_SB1 = all_c_comparison_SB1 %>%
  group_by(SQ_approach) %>%
  summarise(
    mean_error_c = mean(error),
    median_error_c = median(error),
    max_error_c = max(error),
    quantile_error_c = quantile(error, probs = c(0.95)),
    n_dp = length(error),
    n_dp_less_than_ten = length(error[error<10 & error>1]),
    percentage_less_than_ten = (length(error[error<10 &
error>1]))/(length(error)))%>%
  mutate(sample = "SB1") %>%
  ungroup()

```