

RSC Advances



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. This *Accepted Manuscript* will be replaced by the edited, formatted and paginated article as soon as this is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

Correlation of Selected Molecular Properties and Recovery Values in Volatile Organic Compounds Analysis: Comparison of Two Water Matrices

Ivana Ivančev-Tumbas¹, Tatjana-Djaković Sekulić¹, Jelena Molnar, Aleksandra Tubić, Jasmina Agbaba, Jelena Tričković, Marijana Kragulj

University of Novi Sad Faculty of Sciences, Department for Chemistry, Biochemistry and Environmental Protection, Trg D. Obradovića 3, 21000 Novi Sad, Republic of Serbia

Abstract

This study investigates if certain molecular properties can influence the recovery of 18 volatile organic compounds (VOCs) in water under the applied analytical conditions, a purge and trap gas chromatographic method with mass spectrometric detection (P&T GC MS). Statistical and quantitative structure property relationship (QSPR) analyses were applied to find correlations between molecular parameters and analytical recoveries in two different water matrices (clean water (CW) produced in the laboratory and natural groundwater (GW)) at two different concentration levels. At 1 µg/L, most compounds had higher recoveries in CW than in GW, whereas at 15 µg/L, the recoveries were higher in the GW matrix. Polarity number and hydrophilic factor were correlated with the recovery differences at both concentration levels in GW. Polarity was significant in the distinction of recovery differences for CW and GW matrices at low concentration, while air diffusivity had an acceptable correlation with recovery differences for both matrices at the higher concentration. Further correlation of the recoveries themselves with the molecular properties was made by multivariate linear regression (MLR) resulting in a QSPR model. This was only possible for GW at the low 1 µg/L concentration. Partial least square analysis indicated that hydrophilic factor, polarity, and molecular weight were the most important properties investigated. No significant correlation was found in CW matrix or in the higher concentration level in GW matrix, which implies that the most significant properties might only be relevant for the VOCs recovery at low concentrations and only if the matrix contains other water constituents (DOC, salt).

Keywords: VOC; Purge & Trap-GC/MS; recovery; water matrix; stepwise multiple linear regression; validation

1. Introduction

Volatile organic compounds (VOCs) may be present in different types of water. In surface water they mostly originate from traffic pollution (benzene, toluene, ethylbenzene and xylene, e.g. BTEX compounds) and agricultural or industrial pollution (i.e., chlorobenzenes, 1,2-dichloroethane, etc.), whereas in drinking water they are mostly related to disinfection practices (i.e., trihalomethanes). In groundwater, their presence is usually related to pollution caused by waste leakage or accidental pollution. The importance of their presence is underlined in several publications.¹⁻⁴ These compounds are widely used as solvents, cosmetic products, fuels, furnishings, etc. They can easily end up in wastewaters and may pose a risk to the environment. Some VOCs are regulated by the Water Framework Directive⁵, but not all⁶. A year long survey in 1998/99 by Nikolau et al.⁷ showed that in rivers and lake waters in Greece, VOCs concentrations were usually of the order of several µg/L or lower. In wastewaters their concentrations were higher for some compounds (reaching tens or even hundreds of µg/L).

Nowadays, VOC analysis in water is routine. Gas chromatography is mostly applied in combination with different extraction techniques: head space, purge/trap and solid phase micro-extraction. Different types of detectors can be used, such as flame ionisation, electron capture or mass spectrometric detection.⁸ The technique used in this research - purge and trap gas chromatographic method with mass spectrometric detection

¹ Corresponding authors:

ivana.ivancev-tumbas@dh.uns.ac.rs, Tel: +381 21 485 27 46

tatjana.djakovic-sekulic@dh.uns.ac.rs, Tel. +381 21 485 27 42

University of Novi Sad Faculty of Sciences, Department for Chemistry, Biochemistry and Environmental Protection, Trg D. Obradovića 3, 21000 Novi Sad, Republic of Serbia; Fax: +381 21 454 065;

(P&T GC/MS), is recognised as highly sensitive and reliable for very low concentration ranges of various compounds.⁶

A survey of the available literature shows an abundance of data on the precision and accuracy of the different methods. However, almost no attempts have been made to compare recovery values for different matrices. One exception is a paper by Barco-Bonilla et al.⁴ By comparing calibration curve ratios for different wastewater effluents and liquid chromatography/mass spectrometry grade (LC/MS) water (ratios were within the range of 0.8-1.2 for water spiked with 1 and 5 µg/L), they showed that there were no matrix effects. For this purpose, purge and trap (P&T) was coupled with gas chromatography (GC) triple quadrupole mass spectrometry and was used for the analysis of a membrane bioreactor, extended aeration, maturation pond and anaerobic pond. In Standard Methods for the Examination of Water and Wastewater for Volatile Organic Compounds (Method 6200) published in 1998⁹ one can find data related to recovery values at a concentration level of 0.5 µg/L in reagent water. They are within the range of 85-110% for the compounds-of-interest, with relative standard deviations (RSD) in the range 6-10%. For reagent and raw water, recovery values can be found in EPA 524.2.¹⁰ (higher values for raw water for almost all analytes). Concentration levels are different and depend on the water type.

Comparison of recovery values for different matrices is important from the point of view of analytical efficiency. Changes in the recovery due to the matrix are influenced by water constituents and are not a means for controlling the accuracy and precision of a method. Recovery is the amount of a compound that reaches the GC for analysis relative to the amount that was originally present in the sample. The ideal recovery value is 100%. However, acceptable values range from 70-130% for most analytes in gas chromatography methods. In P&T analysis it is essential to vaporise the substances and partition them into the gas phase. The vapour pressure, solubility and extraction temperature affect the procedure. Once the procedure is established, we assume that the influence of the water matrix, if present, can be understood as a possibility to alter the transfer of the analyte into the gaseous phase due to interactions with the water constituents, while other analytical conditions are kept constant as defined by the method. The aim of this study was to investigate if certain molecular properties could influence recovery in the analysis of 18 volatile organic compounds (VOC) in two different water matrices (i.e., clean water (CW) produced in the laboratory, and natural groundwater (GW)) at two different concentration levels (1 and 15 µg/L). Statistical analysis followed by quantitative structure property relationship analysis (QSPR) was applied. In general, QSPR models present relationships between properties of series of molecules and their structural characteristics (derived either from experiments or theoretically). Correlation of the selected molecular parameters and the recovery (assumed as a property under the given analytical conditions) was not studied with the aim to predict recovery but rather to find out if any of the selected molecular parameters could influence recovery under the applied analytical conditions.

2. Materials and methods

2.1. Water matrices

The clean water matrix (CW) was produced in the laboratory by LABCONCO (WaterPro RO/PS Station, Kansas City, USA) system (water of ASTM (American Society for Testing and Materials) Type I quality). It was then boiled for 15 min and purged for 1 hour by nitrogen. It was stored in a glass bottle and prepared daily. Dissolved organic carbon (DOC) content was less than the practical quantitation limit of 0.5 mg/L, and electrical conductivity was 0.055 µS/cm. Measurements of DOC were performed after sample acidification with concentrated hydrochloric acid to pH=2, according to method SRPS ISO 8245:2007.¹¹ pH measurement was carried out on a WTW InoLab (Weilheim in Oberbayern, Germany) portable instrument.

The natural groundwater (GW) matrix was taken from the Danube riparian area (25-30 m depth) free from VOC. Its characteristics were: 3.6 mg/L DOC, pH 7.02, electrical conductivity 770±101 µS/cm, dry residue of 415 mg/L (at 105°C, according to standard method 2540 B⁹).

2.2. Chemicals and reagents

Standards with concentrations of 2000 µg/mL in methanol as Volatile Organic Compounds Mix 7 (Dr. Ehrenstorfer GmbH, Germany) (chloroform, 1,1,1-trichloroethane (1,1,1-TCE), 1,2-dichloroethane (1,2-DCE), benzene, trichloroethene, bromodichloromethane (BDCM), dibromochloromethane (DBCM), bromoform, 1,4-dichlorobenzene), Volatile Organic Compounds Mix 8 (Supelco) (chlorobenzene, ethylbenzene, o-xylene, m-xylene, p-xylene, tetrachloroethene, toluene, 1,2-dichlorobenzene), vinyl chloride (Supelco) and internal

109 standard (IS) fluorobenzene (Supelco) were used to prepare working solutions. Dilution was done with 99.9%
110 methanol (J.T. Baker, Avantor Performance Materials B.V., Deventer, The Netherlands).

111

112

2.3. Purge and trap conditions

113

114

115

116

117

A Tekmar Dohrmann 3100 Sample Concentrator with Vocarb 3000 trap (Carbopack B, 10 cm/Carboxen-1000, 6 cm/Carboxen 1001, 1 cm, *Supelco*, Sigma-Aldrich Co., St. Louis, USA) was used. 5mL aliquots of sample were dispensed into the 5 mL purging device with a gas tight syringe. The sample was purged with a stream of helium at 37.4 mL/min for 11 min at ambient temperature.

118

119

120

Before each sample analysis, the purge and trap system was baked (270°C for 3 min). This was followed by blank analysis of CW. After sample loading, desorption by heating the Vocarb 3000 trap was carried out at 250°C for 2 min. The injector was set to split mode (30:1).

121

122

2.4. Chromatographic and MSD conditions

123

124

125

126

127

128

129

130

131

132

133

134

Table 1. Target ions and qualifiers in GC/MS analysis of VOCs

Compounds	Target ion	Qualifier ions
Vinyl chloride	62	64
Chloroform	83	85
1,1,1-TCE	97	99
1,2-DCE	62	98
Benzene	78	77
Trichloroethene	95	130, 132
BDCM	83	127, 129
Toluene	91	92
DBCM	129	208, 173
Tetrachloroethene	166	168, 129
Chlorobenzene	112	77
Ethylbenzene	91	106
m+p- xylene	106	91
o- xylene	106	91
Bromoform	173	171, 252
1,2-dichlorobenzene	146	148, 252
1,4-dichlorobenzene	146	148, 252
Fluorobenzene (IS)	96	97

135

136

137

138

2.5. Calibration and quantification

139

140

141

142

143

144

Calibration curves were obtained by spiking CW matrix with 17 VOCs mixture in methanol by gas-tight syringe, in accordance with the instructions given in Standard Methods 6200B⁹, with further sample processing of the samples as explained in Standard Methods 6200B⁹. 9 point calibration curves were made using concentrations of 0.4, 0.5, 1.0, 2.0, 5.0, 8.0, 10, 14 and 25 µg/L and for vinyl chloride within the range of 0.2-25 µg/L. The internal standard concentration was 10 µg/L. For the curves in the 0.4-25 µg/L concentration range, the coefficients of determination (R^2) were within the range of 0.991-0.998.

145

146

147

148

As a part of routine quality control, a calibration verification standard (CVS) from the same source and a laboratory control standard (LCS) from an independent source were used for verification. A criterion of $\pm 15\%$ of difference in comparison to the initial calibration was accepted in accordance with EPA 8000B.¹² The obtained results indicated that all the measured values were within the range of $\pm 15\%$ of the expected value.

149 RSDs for 50 measurements of the calibration verification standard (4 $\mu\text{g/L}$ CVS) were collected and for all
150 compounds were lower than 10%. LCS measurements ranged from 70% to 130% of the expected concentration,
151 in accordance with EPA 8000B.¹²

152
153
154

155 2.6. Method performance

156
157
158
159
160
161
162
163

EPA Methods 5030B¹³ and 8260B¹⁴ were used to develop an internal laboratory procedure for the analysis. Method detection limits were determined along the guidelines given by Glase et al.¹⁵ (analyte added in a concentration which is 1-5 times the estimated method detection level (MDL)). The MDL was determined based on the RSD of 6 measurements of both spiked CW and GW matrix at a concentration level of 0.4 $\mu\text{g/L}$, except for vinyl chloride (0.2 $\mu\text{g/L}$). The practical limit of quantitation (PQL) was calculated as 5 x MDL according to EPA Method 8260B.¹⁴ Repeatability for this concentration level was determined as a RSD of 6 consecutive measurements.

164
165
166
167
168

Precision was assessed as the RSD of recovery values determined for four samples analysed in 6 series in duplicates: clean water matrix at concentrations of 1 $\mu\text{g/L}$ and 15 $\mu\text{g/L}$ (CW_1 and CW_{15}) and groundwater matrix at the same concentrations (GW_1 and GW_{15}). One series consisted of a blank sample (non-spiked CW) and duplicates of the four mentioned samples. Additionally, the groundwater matrix (GW) was checked for the presence of VOC.

169
170
171

172 2.7. Recovery comparison

173
174
175
176

The recovery comparison for the two water matrices was carried out by statistical analysis of 18 VOCs for 12 replicates at both concentration levels, i.e. 1 and 15 $\mu\text{g/L}$. The concentration levels selected are similar to the VOC levels defined as relevant environmental quality standards for surface waters (Directive 2008/105/EC)¹⁶.

177
178
179
180
181
182
183
184
185

Nine molecular parameters (descriptors) (water solubility (S, mg/L), Henry's law constant (HLC, $\text{atm}\cdot\text{m}^3/\text{mol}$), octanol-water partition coefficient (K_{ow}), air ($D_{i,a}$) and water diffusivity at 25°C ($D_{i,w}$) taken from USEPA,¹⁷ as well as molecular mass (MW), polarity number (Pol), hydrophilic factor (Hy) and molar refractivity calculated by Dragon¹⁸) were selected for the study. Descriptor values are given in Supplementary Table S2. The influence of molecular specific parameters on differences in recovery, i.e. its median value, was analysed by finding correlations between differences in recovery for the following: two concentrations in clean water (CW_{1-15}), two concentrations in groundwater (GW_{1-15}), a lower concentration in two different matrices (CW-GW_1) and a higher concentration in two different matrices (GW-CW_{15}) for compounds 1-18. In addition, the recovery values themselves were correlated with the molecular descriptors.

186
187
188
189
190
191
192
193
194

In order to visualize similarities and differences in water matrices, hierarchical clustering analysis (HCA) was used (Ward's linkage method). Furthermore, linear regression followed by stepwise multiple linear regression (MLR) (forward selection method) was applied^{19,20}. QSPR models for median recovery values based on a few deliberately selected explanatory variables (molecular descriptors) were constructed. The data were organized in matrices \mathbf{X} (18x9) where the rows represented the 18 investigated compounds (VOC), and the columns corresponded to the 9 molecular parameters. The independent variables were the 9 selected descriptors while the dependent variables were the observed median values of recovery for 12 measurements of each of the 18 VOCs. Additionally, the partial least squares (PLS) method was used since it is a well-known approach for the analysis of multidimensional data sets.²¹⁻²⁵

195
196
197
198

All the calculations were carried out by STATISTICA v. 10.0.²⁷ The data were mean-centred (subtracting the mean and dividing by the standard deviation²⁶) before any statistical operation in order to prevent the highly abundant components from dominating in the final result over the components present in much smaller quantities.

199
200
201

201 3. Results and discussion

202
203
204

203 3.1. Method performance

205
206
207

All the recovery values at the 0.4 $\mu\text{g/L}$ concentration level (0.2 $\mu\text{g/L}$ for vinyl chloride) were within the acceptable range (70-130%), with the exceptions of chloroform and toluene, due to the laboratory contamination. This contamination affected only the low 0.4 $\mu\text{g/L}$ concentration level. This caused the

208 determined MDL and PQL values for these analytes to be slightly higher (i.e. PQLs of 1.60 and 1.06 $\mu\text{g/L}$
209 respectively). Generally speaking, the calculated MDL and PQL values were within the range found in Standard
210 Methods⁹ and EPA Method 8260B¹⁴. In the case of BTEX, chlorobenzenes, chloroform, toluene,
211 trichloroethylene and tetrachloroethylene, the values were somewhat higher. Details on the method performance
212 and corresponding values from the literature are given in supplementary material S1. Based on the calculation
213 procedure explained in 2.6, different MDLs were obtained for different matrices. The MDL values in CW were
214 within the range of 0.027-0.32 $\mu\text{g/L}$, while for GW the range was somewhat narrower (0.026-0.18 $\mu\text{g/L}$), and
215 for most of the compounds the MDLs were lower (except vinyl chloride, 1,2-DCE and 1,1,1-TCE where there
216 were no differences between the samples), whereas for benzene it was higher. Repeatabilities for the 0.4 $\mu\text{g/L}$
217 concentration level (0.2 $\mu\text{g/L}$ for vinyl chloride) were determined as the RSD of 6 consecutive measurements
218 and were within the range of 1.2-13.2%.

219 It should be noted that the values from table S1 (given for evaluation of method performance) were not the
220 basis for further recovery comparison at concentration levels of 1 and 15 $\mu\text{g/L}$. The recovery values for
221 concentrations of 1 and 15 $\mu\text{g/L}$ in different matrices were collected and further evaluated by QSPR. Average
222 values for 12 replicates of each sample were in the range of 83-123% for CW matrix and in the range of 72-
223 108% in GW matrix. Method precision at these concentrations was within the range of RSD 1.0-8.7% for CW
224 matrix and within the range of RSD 2.0-9.4% for GW matrix.

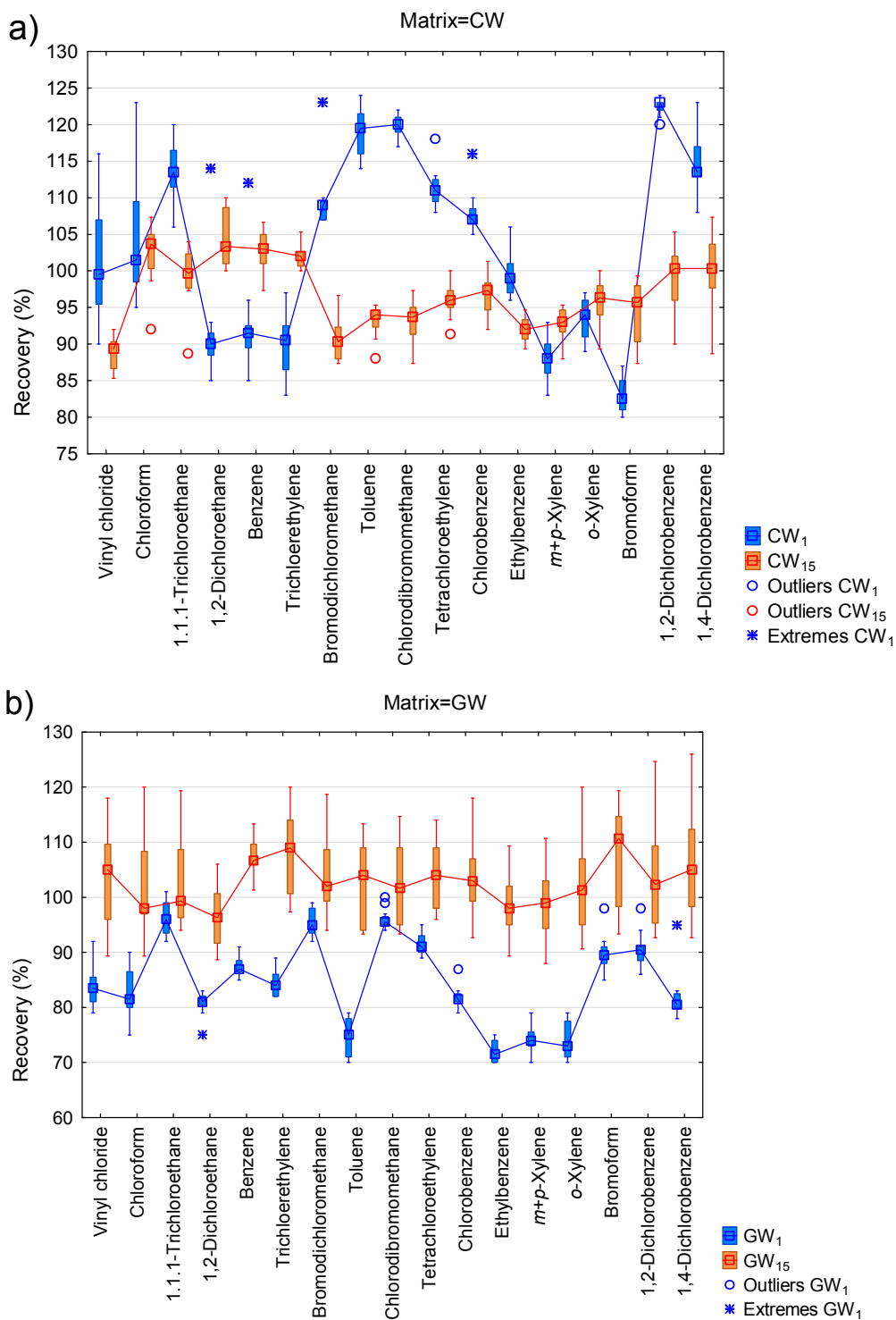
225
226
227

228 3.2. Comparison of recoveries

229

230 Box plots of recovery values for the 12 replicates of each substance at concentration levels of 1 and 15 $\mu\text{g/L}$
231 and both matrices grouped by compound are given in Figure 1 and Figure 2. A parallel presentation of all the
232 performed measurements allows comparison of recovery values of individual samples. Boxes themselves and
233 their width indicate the degree of dispersion and skewness in the data. Central tendency (median) is connected
234 by line in order to facilitate comparison of the results obtained in different experimental conditions. In addition,
235 the quartiles are given as an indication of the dispersion of the recovery values. More precisely, the bottom and
236 the top of the box are the lower (Q_1) and upper (Q_3) quartiles, respectively. The ends of the whiskers represent
237 1.5 interquartile ranges ($\text{IQR} = Q_3 - Q_1$). The values below the lower whiskers ($Q_1 - 1.5 \text{ IQR}$) and above the upper
238 whiskers ($Q_3 + 1.5 \text{ IQR}$) are outliers. Consequently, extreme outliers are values beyond 3 IQR from Q_1 and Q_3 .
239 The asymmetrical position of the median in the boxes in Figures 1 and 2 is an indication that the recovery data
240 do not follow a normal distribution (most probably due to the limited number of data points) and hence
241 convenient descriptive statistics, such as mean and standard deviation, are not appropriate.

242



243

244

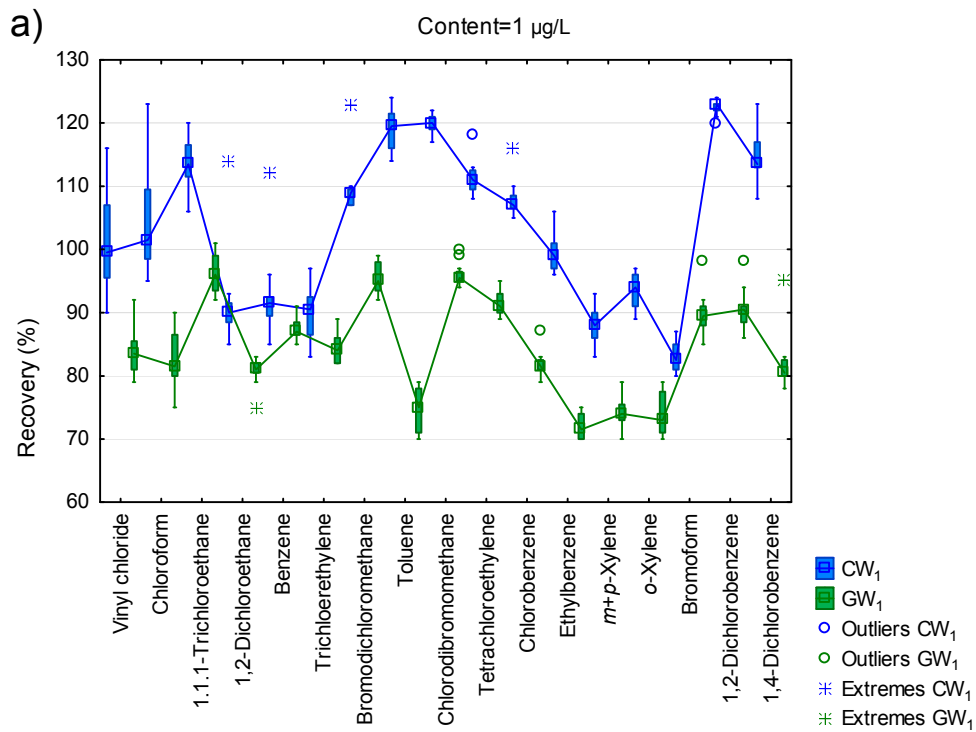
245

246

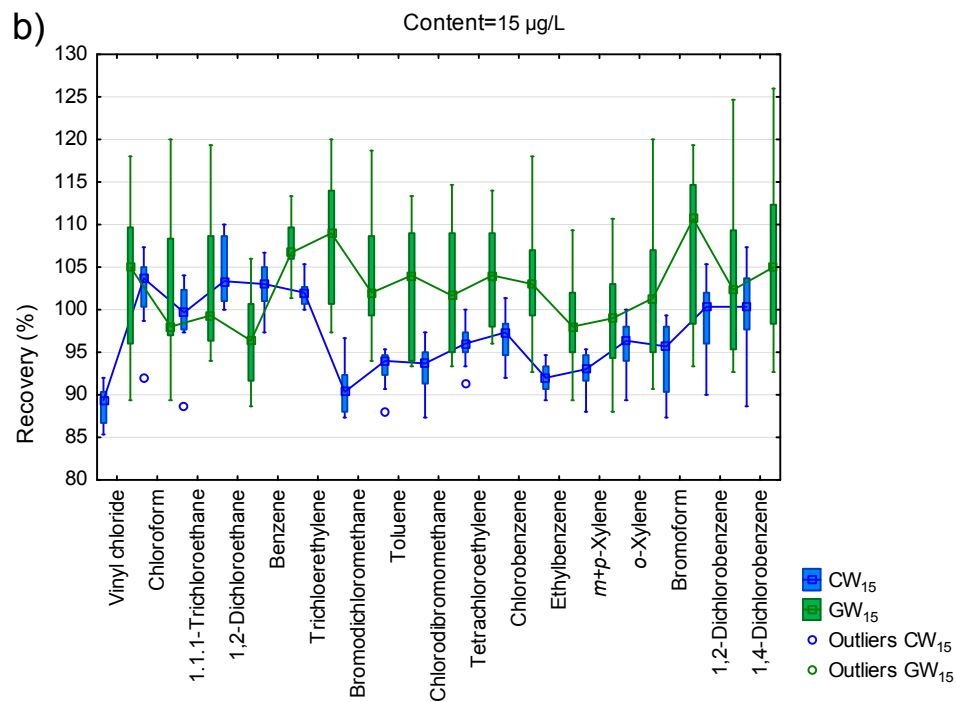
247

248

Fig 1. Box plot of recovery data for 12 replicates of 18 VOCs for clean water matrix at concentration levels of 1 $\mu\text{g/L}$ (CW₁) and 15 $\mu\text{g/L}$ (CW₁₅) (a) and for groundwater matrix at concentration levels of 1 $\mu\text{g/L}$ (GW₁) and 15 $\mu\text{g/L}$ (GW₁₅) (b).



249
250



251
252
253
254
255
256
257
258

Fig 2. Box plot of recovery data for 12 replicates of 18 VOCs for for both clean and groundwater matrices at the concentration level of 1 $\mu\text{g/L}$ (a) and 15 $\mu\text{g/L}$ (b).

259 In clean water matrix (Figure 1a) for 1,2-dichloroethane, benzene, trichloroethylene and bromoform,
260 recoveries are higher at the higher concentration level, while for vinyl chloride, 1,1,1-trichloroethane,
261 bromodichloromethane, toluene, chlorodibromomethane, tetrachloroethylene, chlorobenzene, ethylbenzene, 1,2-
262 dichlorobenzene and 1,4-dichlorobenzene, they are lower at the higher concentration. In ground water matrix
263 (Figure 1b), all the compounds have higher recoveries at the higher concentration level.

264

265 Comparison of recovery values at the low concentration level for different matrices (Figure 2a) showed that
266 the significant differences were not observed for benzene, trichloroethylene and bromoform, while for all the
267 other compounds the values were higher in CW than in GW. In contrast, at the higher concentration level
268 (Figure 2b), recoveries were higher in GW.

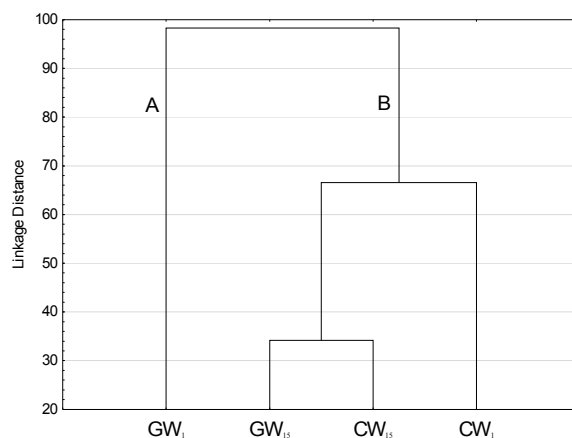
269

270 In all further consideration the medians were used, since they are not influenced by extreme outliers in the
271 data set, which is not the case with arithmetic mean values.

272

273 HCA was performed in order to visualize similarities (clusters formation) and differences (linkage distances
274 between clusters) of the recovery values for the two matrices and two concentration levels. The resulting
275 dendrogram is presented in Figure 3. The y-axis represents the corresponding linkage distances between the two
276 objects or clusters which are merged. The dendrogram reveals two distinct clusters: cluster A with recoveries
277 related to GW_1 and cluster B that included recoveries for CW_1 together with one sub-cluster with recoveries
278 observed in matrices spiked with higher concentration levels (GW_{15} and CW_{15}). One can hypothesize that the
279 recovery values for GW_1 might be influenced by other factors than in the other three cases since they are
280 grouped in two separate clusters- A and B. Furthermore, the difference between the two matrices at the higher
281 concentration level is lower than at the lower concentration level of $1 \mu\text{g/L}$.

282



283

284

285 **Fig 3.** Dendrogram of two water matrices at two concentration levels (GW_1 - groundwater matrix with the $1 \mu\text{g/L}$
286 concentration level, GW_{15} - groundwater matrix with the $15 \mu\text{g/L}$ concentration level, CW_1 - clean water matrix
287 with the $1 \mu\text{g/L}$ concentration level, CW_{15} - clean water matrix with the $15 \mu\text{g/L}$ concentration level) in the space
288 of 18 VOC compounds obtained by the Ward linkage method using Euclidean distance

289

3.2.1. Influence of chemical-specific molecular parameters on recovery

290

291 The differences evident in the recovery values obtained for the two matrices as well as the two concentration
292 levels may result from certain molecular characteristics that define the general behaviour of the compounds
293 during analysis, assuming that all the other analytical conditions are held constant and defined by the method.
294 These different intrinsic properties of molecules are usually expressed by various molecular descriptors. Their
295 selection is governed by general knowledge of the features that influence the transfer from water to air and
296 hence have a decisive influence on analytical recoveries while all the other analytical conditions are kept
297 constant. For example, transfer from water to air might be influenced by interactions between analytes and water
298 constituents. Our assumption is that water DOC is not extracted, nor purged nor injected into the GC during
299 P&T GC/MS analysis, since it is known from the literature that in natural groundwater a negligible portion of
300 the DOC is usually volatile. This *a priori* excludes the kind of influence by sorptive interaction of DOC with the
301 relevant instrument surfaces such that the only possible interactions are those of solutes and constituents in the

302 water phase. Correlations were therefore used as an exploratory tool to indirectly assess if the selected molecular
303 properties might influence analytical recoveries under controlled analytical conditions.

304 The influence of chemical molecular descriptors on recovery values was investigated for different matrices,
305 as well as for the different concentration levels for each matrix. Initially, the relationship between the
306 differences in recovery for the different concentrations in clean water (CW_{1-15}), for the different concentrations
307 in groundwater (GW_{1-15}), for the lower concentration in the two different matrices ($CW-GW$)₁ and for the higher
308 concentration in the two different matrices ($GW-CW$)₁₅ for VOCs and the selected chemical-specific parameters
309 was examined. The intention was to find out which parameter describes the observed differences with the best
310 fit. According to Pearson's correlation (Table 2), linear correlations for the observed differences in recovery and
311 molecular descriptors appear in a wide range, but not higher than 0.68. According to Bevington²⁸ the limiting
312 value for linear-correlation coefficients for an acceptable correlation of the 18 compounds at a probability level
313 of 0.05 is a minimum 0.468. Correlations exceeding this limit are bolded in Table 2 and are significant (high
314 enough for acceptable linear correlation). Compared to other molecular descriptors, polarity number (Pol) and
315 hydrophilic factor (Hy) had a significant linear correlation coefficient for the difference observed at the two
316 concentration levels in GW matrix. Pol is significant for the difference observed between CW and GW matrix at
317 the lower concentration level (bold in Table 2). At the higher concentration level, differences between the CW
318 and GW matrix showed significant linear correlation with air diffusivity ($D_{i,a}$).
319

320 **Table 2.** Significance of molecular specific parameters for recovery differences between two matrices and two
321 concentration levels expressed by Pearson's correlation ($p < 0.05$, $N=18$).
322

Recovery difference	S (mg/L)	log Kow	$D_{i,a}$ (cm ² /s)	$D_{i,w}$ (cm ² /s)	MW	Pol	Hy	molar refractivity	HLC (atm-m ³ /mol)
CW_{1-15}	0.2251	-0.1768	-0.2485	0.1754	-0.1490	-0.0672	-0.0763	-0.1065	-0.1059
GW_{1-15}	-0.4063	0.3954	0.3499	-0.2330	-0.4185	0.5713	-0.6843	0.3802	-0.0030
($CW-GW$) ₁	-0.3489	0.4364	0.1945	-0.1774	-0.2437	0.4831	-0.3691	0.3545	-0.0209
($GW-CW$) ₁₅	0.3603	-0.0589	-0.4974	0.3589	-0.2752	-0.0532	-0.0053	-0.0408	0.2821

323 S- water solubility, mg/L; HLC- Henry's law constant, atm-m³/mol; K_{ow} - octanol-water partition coefficient;
324 $D_{i,a}$ - air and $D_{i,w}$ - water diffusivity at 25°C taken from USEPA,¹⁷ MW- molecular mass; Pol- polarity number;
325 Hy- hydrophilic factor
326

327 Subsequently, an attempt was made to correlate not the differences in recoveries for certain samples, but the
328 recovery values themselves and more than one molecular descriptor by MLR. The goal was not to predict the
329 recovery values but to investigate which among the selected molecular descriptors might influence them and
330 under what conditions.
331

332 Statistical parameters for the models calculated by stepwise regression on standardized data are shown in Table
333 3. The developed models were internally validated. The cross-validated regression coefficient (Q^2) is defined as:
334

$$335 Q^2 = 1 - \frac{\sum (Y_{pred} - Y_{exp})^2}{\sum (Y_{exp} - Y_{mean})^2} \quad (1)$$

336 where Y_{pred} , Y_{exp} and Y_{mean} are the predicted, experimental, and mean values of the target property
337 (recovery).
338

339 Q^2 higher than 0.5 is a necessary condition for the model to have predictive power, but still does not
340 automatically imply high predictability.²⁹ Additional important cross-validation parameters accounting for a
341 good estimate of the real predictive error of the model are PRESS (predictive residual error sum of squares),
342 SSY (sum of squares of deviation of the dependent variable values from their mean) and PRESS/SSY that
343 should be smaller than 0.4.³⁰ In Table 3, Model 1 shows the best statistics, although it is a borderline case for the
344 PRESS/SSY ratio. Model 3 failed because it does not fulfil the statistical criteria ($Q^2 > 0.5$). The performance of
345 model 4 cannot be accepted as it is over-parameterized (the ratio n/descriptors is >5).
346
347

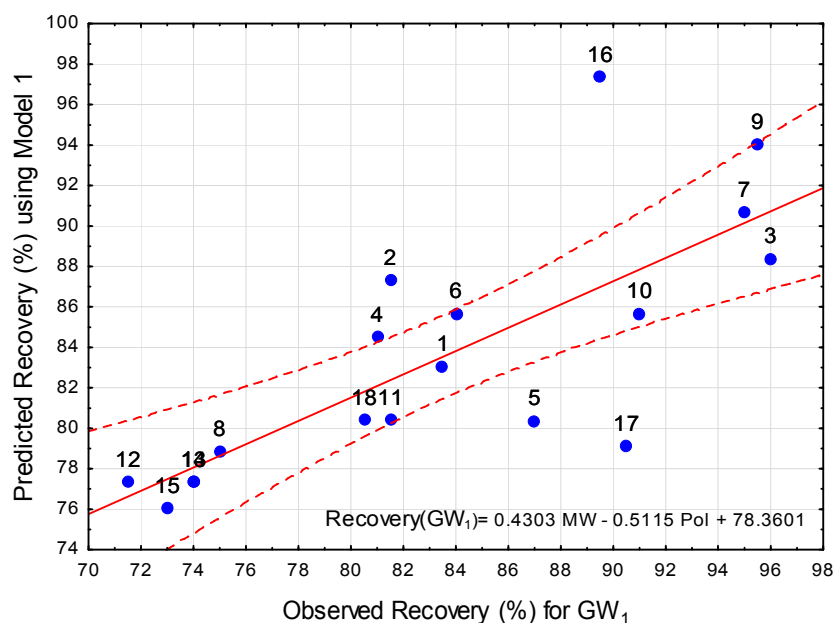
348 **Table 3 .** Internal validation statistics¹ for obtained QSPR models

	Molecular descriptor included in model	Multiple - R	Q^2	Adjusted - R^2	F	p	PRESS	PRESS/SSY	Model No.
GW_1	MW, Pol	0.7589	0.5760	0.5194	10.1869	0.0016	473.4190	0.424042	1
	$D_{i,a}$, Pol	0.7484	0.5601	0.5014	9.5481	0.0021	491.1586	0.439931	2

GW ₁₅	S, Pol	0.5721	0.3274	0.2377	3.6500	0.0511	174.1626	0.672648	3
CW ₁₅	S, HLC, D _{i, as} , MW, Pol	0.8893	0.7909	0.7038	9.0776	0.0009	74.8705	0.289164	4

349 ¹R-correlation coefficients, R^2_{adj} - adjusted square of the correlation coefficient, p -significance level (set to
 350 $p \leq 0.05$), and Fisher test for significance of the equation (F-test). The acceptance level for the individual
 351 independent variable was set to 95% significance level.
 352

353 In Figure 4 the best constructed stepwise regression model describing recovery for the tested compounds in
 354 GW₁ (model 1) is presented.
 355
 356



357
 358 **Fig 4.** Observed vs. predicted values of Recovery according to QSPR Model 1 from Table 3 (1 – vinyl chloride;
 359 2 – chloroform; 3 – 1,1,1-trichloroethane; 4 - 1,2-dichloroethane; 5 – benzene; 6 – toluene; 7 –
 360 bromodichloromethane; 8 – tetrachloroethylene; 9– chlorodibromomethane; 10 – trichloroethylene; 11 –
 361 chlorobenzene; 12 – ethylbenzene; 13 - *m*-xylene; 14 - *p*-xylene; 15 - *o*-Xylene; 16 – bromoform; 17 - 1,2-
 362 dichlorobenzene; 18 - 1,4-dichlorobenzene) with indicated 95% confidence interval band to the regression line.
 363

364 Recovery values for compounds 2 (chloroform), 5 (benzene), 16 (bromoform) and 17, (1,2-dichlorbenzene) are
 365 not well described by the model proposed (see Figure 3). Sorting the values for each descriptor in increasing
 366 order one could find that these substances may be grouped only by HLC (all of them belong to the subgroup
 367 which has values from $5.4 \cdot 10^{-4}$ to $5.6 \cdot 10^{-3}$ atm·m³/mol). However, six more substances belong to that group (*o*-
 368 xylene, chlorobenzene, 1,4-dichlorobenzene, bromodichloromethane, 1,2-dichloroethane and
 369 chlorodibromomethane, numbered 15, 11, 18, 7, 4 and 9 respectively) but fit well with the proposed model. The
 370 rest of the compounds have higher HLC values up to $2.7 \cdot 10^{-2}$ atm·m³/mol.
 371

372 MLR can be a good method for data analysis in cases when the descriptors are few in number, not significantly
 373 collinear and if the equation has a good fit. However, if the number of descriptors is too large (in our case more
 374 than 3 for 18 substances, like in model 4) MLR becomes inappropriate. Therefore we applied PLS regression.
 375 The difference between MLR and PLS is that PLS fits the model simultaneously for all descriptors, while MLR
 376 fits descriptors separately.³¹ On the other hand, comparison between MLR and PLS models is more abstract
 377 which makes it difficult to understand and interpret. PLS regression is an extension of the multiple linear
 378 regression model. The main purpose is to build a linear model transforming the original variables into the new
 379 orthogonal variables, maximizing the description of a covariance between X and Y²⁶.
 380

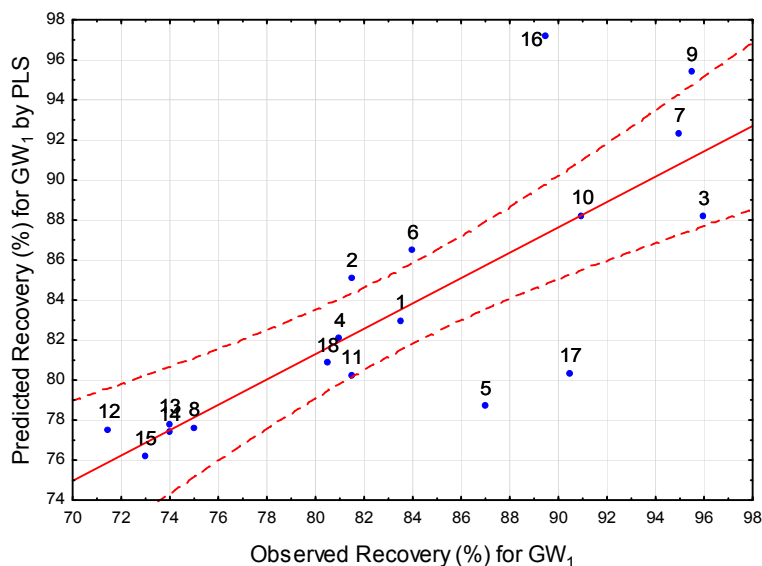
381 It was only possible to construct a PLS model for GW₁, which is in accordance with findings related to results
 382 presented in Figure 3. Statistical data for the PLS model are given in Table 4. 63.33% of the sum of squares of
 383 the dependent variables is explained by two PLS components. Figure 5 shows the observed vs the predicted

384 recovery values according to the PLS regression model for the tested compounds in GW₁. Compared to Figure
 385 4, Figure 5 does not seem very much improved particularly for compounds 5 (benzene), 9
 386 (chlorodibromomethane), 16 (bromoform) and 17 (1,2-dichlorobenzene). Improvement was achieved for
 387 compounds 4 (1,2-dichloroethane) and 10 (trichloroethylene).
 388

389 **Table 4.** Statistical data for PLS model related to GW₁
 390

	R ² X(Cumul.)	Eigenvalues	R ² Y(Cumul.)	Q ² (Cumul.)	Significance	Iteration
1	0.4254	3.5865	0.5517	0.4080	S	1
2	0.6522	1.6407	0.6333	0.1449	NS	1

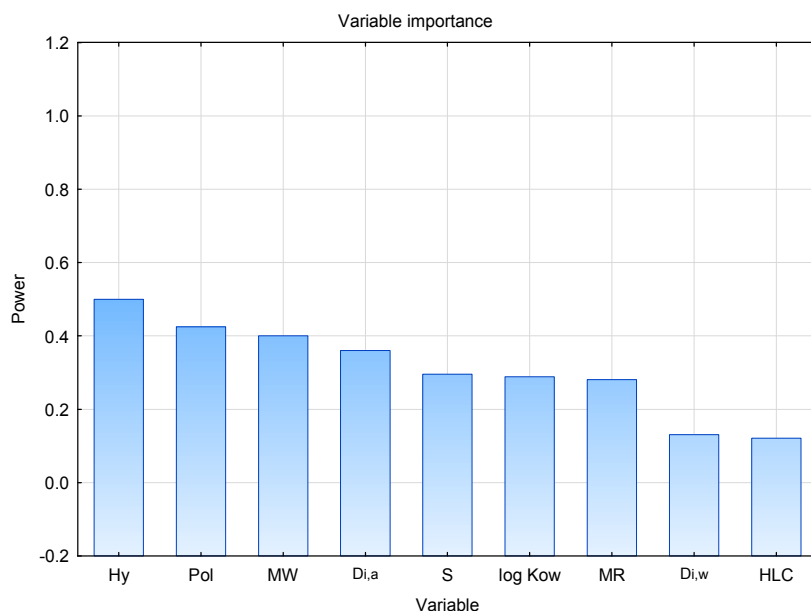
391
 392



393
 394
 395
 396
 397
 398
 399
 400
 401

Fig 5. Observed vs. predicted Recovery values for PLS model obtained for groundwater at concentration level of 1µg/L (GW₁) (1 – vinyl chloride; 2 – chloroform; 3 – 1,1,1-trichloroethane; 4 - 1,2-dichloroethane; 5 – benzene; 6 – toluene; 7 – bromodichloromethane; 8 – tetrachloroethylene; 9– chlorodibromomethane; 10 – trichloroethylene; 11 – chlorobenzene; 12 – ethylbenzene; 13 - *m*-xylene; 14 - *p*-xylene; 15 - *o*-xylene; 16 – bromoform; 17 - 1,2-dichlorobenzene; 18 - 1,4-dichlorobenzene) with indicated 95% confidence interval band to the regression line

402 In order to find out the significance of molecular specific parameters in the PLS model the variable importance
 403 for the projection (VIP) is given in Fig 6. This allows quick identification of the explanatory variables
 404 (descriptors) that contribute most to the model. The variables with higher VIP scores (e.g. power) are considered
 405 more relevant.
 406
 407



408
409 **Fig 6.** The VIP scores for molecular descriptors used in the PLS calculation for water matrix GW1
410

411 Thus the results obtained by PLS calculation show that the hydrophilic factor has the highest VIP score followed
412 by polarity, molecular weight and Di,a which were already indicated by the MLR.
413

414 Based on these results one can conclude that recovery values can only be correlated with some of the selected
415 molecular properties in the case of groundwater which contains DOC, analyte at a low concentration level
416 ($1\mu\text{g/L}$) and salts (GW_1). A possible explanation might be related to the potential for physical or chemical
417 interaction of analytes and water constituents to the extent that could influence the recovery in an indirect way.
418 In the case of CW there is most probably no possibility for such interactions to occur so it is only logical that the
419 selected molecular properties have no influence on the recovery under the analytical conditions applied. As for
420 the groundwater matrix at higher concentrations of analytes (GW_{15}) and the same DOC and salt concentrations
421 as the low concentration level of spiked VOCs (GW_1), one possible reason for the absence of correlation with
422 the selected descriptors and recoveries might be a limited capacity for relevant interactions to alter the transfer
423 of analyte into the gas phase. Thus the influence might become negligible in comparison to GW_1 . Testing this
424 hypothesis would require further work, at more concentration levels and in different natural matrices with
425 varying amounts of DOC amount and other characteristics, which is beyond the scope of this current work.
426

427 428 **4. Conclusions**

429
430 Statistical analysis of the dataset showed that the recovery values for VOCs have different behaviours depending
431 on the matrix type and concentration level. Comparison of the recovery values at the low concentration level
432 ($1\mu\text{g/L}$) showed that for the majority of compounds the recovery values are higher in CW than in GW matrix. In
433 contrast, at the higher concentration level ($15\mu\text{g/L}$), recovery values are higher in GW matrix.
434

435 Using QSPR analysis, polarity number and hydrophilic factor were found to have an acceptable linear
436 correlation with differences between the recovery values at two concentration levels in GW matrix. Polarity was
437 significant in the distinction of recovery differences for CW and GW matrix at the low concentration level,
438 while air diffusivity had an acceptable value in the distinction of recovery differences for the two matrices at the
439 higher concentration.
440

441 Correlation between the molecular properties and the recovery values by MLR resulted in a QSPR model for the
442 groundwater matrix at the low VOC concentration. The QSPR includes molecular weight and substance polarity
443 as molecular parameters. This result was further supported by PLS, although with somewhat different results for
444 potentially important molecular parameters based on VIP scores: hydrophilic factor > polarity > molecular
445 weight. No significant correlation was found in either CW matrix or in GW matrix at the higher concentration

446 level. This implies that these most significant properties of VOC molecules might be relevant for the recovery
447 only at low VOC concentration levels and only if the matrix contains other water constituents (i.e. DOC, salts).
448

449

450 Acknowledgements

451

452 This research is funded by Ministry of Education, Science and Technological Development of Republic of
453 Serbia (grant ON 172028). The authors would like to thank prof. Zagorka Lozanov-Crvenković for valuable
454 comments and constructive suggestions during the preparation of this manuscript.

455

456 References

457

- 458 1. S. Lacorte, L. Olivella, M. Rosell, M. Figueras, A. Ginebreda and D. Barcelo, *Chromatographia*, 2002, **56**,
459 739.
- 460 2. E. Martinez, S. Lacorte, I. Llobet, P. Vianna and D. Barcelo, *J. Chromatogr. A.*, 2002, **959**, 181.
- 461 3. A. Ikem, *J. Food Compos. Anal.*, 2010, **23**, 70.
- 462 4. N. Barco-Bonilla, P. Plaza- Bolaños, J.L. Fernandez-Moreno, R. Romero-Gonzales, A.G. Frenich and J.L.
463 Martinez Vidal, *Anal. Bioanal. Chem.*, 2011, **400**, 3537.
- 464 5. WFD, Water Framework Directive 2000/60/EC, OJ L 327/1, 2000.
- 465 6. M. Rosell, S. Lacorte, A. Ginebreda and D. Barcelo, *J. Chromatogr. A.*, 2003, **995**, 171.
- 466 7. Nikolau A.D., Goufopoulos S.K., Kostopoulou M.N., Kolokythas G.A., Iekkas T.D. (2002) *Water Research*
467 **36**, 2883-2890.
- 468 8. P. Lepom, B. Brown, G. Hanke, R. Loos, P. Quevauviller and J. Wollgast, *J. Chromatogr. A.*, 2009, **1216**,
469 302.
- 470 9. AWWA-APHA-WEF Standard Methods for the Examination of Water and Wastewater, 20th ed., American
471 Public Health Association/American Water Works Association/Water Environment Federation, Washington,
472 DC, 1998. Volatile organic compounds (6200 B), pp 6-25
- 473 10. USEPA Method 524.2 Measurement of pureable organic compounds in water by capillary column gas
474 chromatography/mass spectrometry, Revision 4, 1992.
- 475 11. SRPS ISO 8245:2007 Guidelines for determination of total organic carbon (TOC) and dissolved organic
476 carbon (DOC) in water, 2007.
- 477 12. USEPA Method 8000B Determinative chromatographic separations, Revision 2, 1996b.
- 478 13. USEPA Method 5030B Purge-and-trap for aqueous samples, Revision 2, 1996a.
- 479 14. USEPA EPA Method 8260B Volatile organic compounds by gas chromatography/mass spectrometry
480 (GC/MS), Revision 2, 1996c.
- 481 15. A.A. Glase, D.L. Foerst, G.D. McKee, S.A. Quave and W.L. Budde WL, *Environ. Sci. Tech.*, 1981, **15**, 1426.
- 482 16. Directive 2008/105/EC, OJ L348/84, 24.12.2008
- 483 17. USEPA Supplemental guidance for developing soil screening levels for superfund sites. United States
484 Environmental Protection Agency Solid Waste and Emergency Response, 2002.
485 http://www.epa.gov/superfund/health/conmedia/soil/pdfs/part_5.pdf. Accessed 15 June 2012.
- 486 18. Talete srl DRAGON (Software for Molecular Descriptor Calculation) Version 6.0 - 2011 -
487 <http://www.talete.mi.it>, 2011.
- 488 19. D. Cox, *J. Royal Stat. Society: Series B*, 1972, **34**, 187.
- 489 20. S. Bergante, G. Faccioto and G. Minotta, *Central Europe Journal of Biology*, 2010, **5**, 522.
- 490 21. D. Massart, B. Vandeginste, L. Buydens, S. DeJong, P. Lewi and J. Weyers-Verbeke, *Handbook of*
491 *Chemometrics and Qualimetrics: Part A*, Elsevier: Amsterdam, 1997.
- 492 22. H. Martens and T. Naes, *Multivariate calibration by data compression*. In: Williams PC, Norris K (eds)
493 Near-infrared Technology in Agricultural and Food Industries, American Association of Cereal Chemists: St.
494 Paul, Minnesota, 1987.
- 495 23. H. Martens and T. Naes, *Multivariate Calibration*. John Wiley & Sons: New York, 1989.
- 496 24. H. Wold, *Soft modeling: the basic design and some extensions*. In: Jöreskog K, Wold H (eds) Systems Under
497 Indirect Observation: Causality, Structure, Prediction, North-Holland: Amsterdam, 1981.
- 498 25. S. Wold, H. Martens and H. Wold, *The Multivariate Calibration Problem in Chemistry Solved by the PLS*
499 *Method*. Lecture Notes in Mathematics. Springer Verlag: Heidelberg, 1983.
- 500 26. P. Geladi, B. R. Kowalski, *Analytica Chimica Acta*, Volume 185, 1986, Pages 1-17
- 501 27. StatSoft Inc. STATISTICA (data analysis software system), version 10. www.statsoft.com, 2011, Tulsa, OK
502 USA
- 503 28. P. R. Bevington, *Data Reduction and Error Analysis for the Physical Sciences*, McGraw Hill Book Co.,
504 New York, 1969.

- 505 29. A. Golbraikh and A. Tropsha, *J. Mol. Graph. Model*, 2002, **20**, 269.
506 30. S. Wold, *Mol. Inform.*, 1991, **10**, 191.
507 31. T. Dijkstra, *J. Econometrics*, 1983, **22**, 67.
508