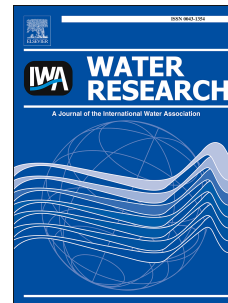


Accepted Manuscript

Predicting chloroform production from organic precursors

Tom Bond, Nigel Graham



PII: S0043-1354(17)30636-X

DOI: [10.1016/j.watres.2017.07.063](https://doi.org/10.1016/j.watres.2017.07.063)

Reference: WR 13106

To appear in: *Water Research*

Received Date: 3 May 2017

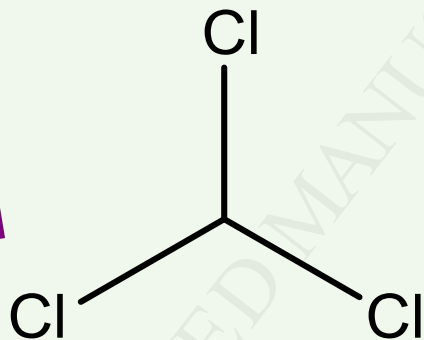
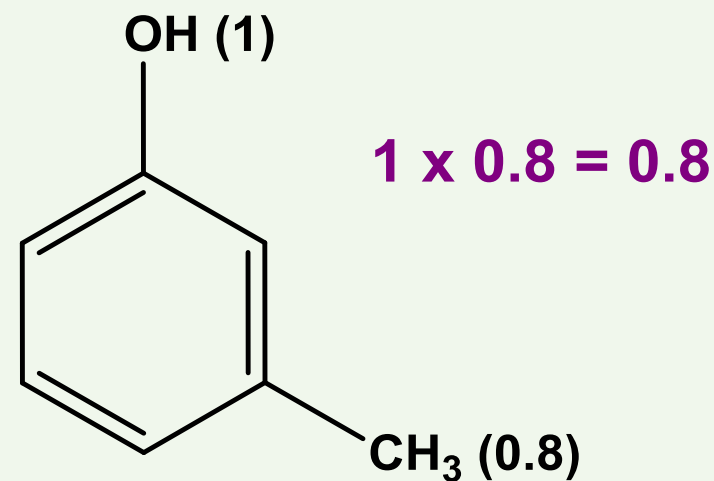
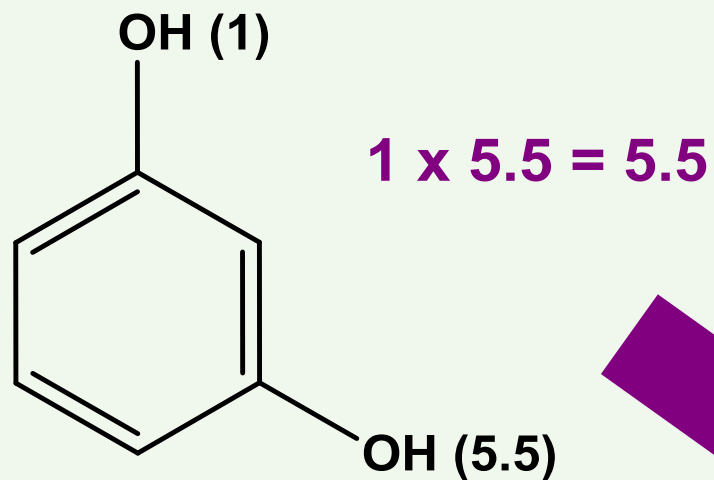
Revised Date: 17 July 2017

Accepted Date: 23 July 2017

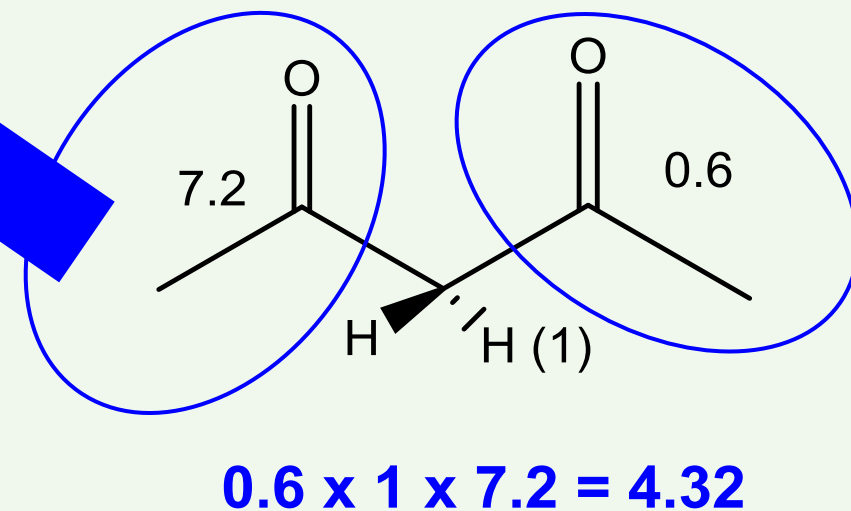
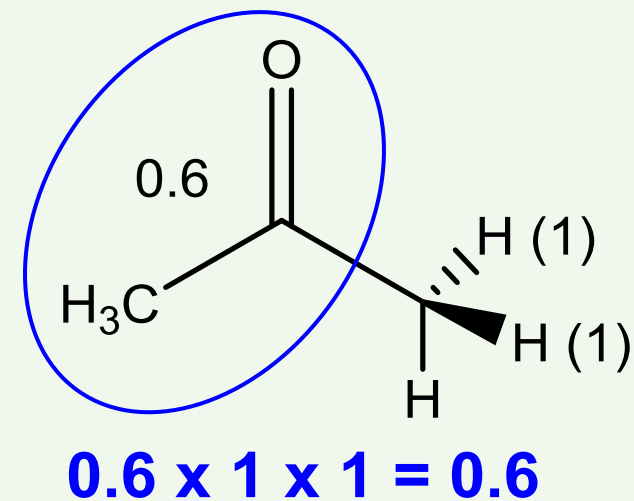
Please cite this article as: Bond, T., Graham, N., Predicting chloroform production from organic precursors, *Water Research* (2017), doi: 10.1016/j.watres.2017.07.063.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Aromatic score



Enolizable score



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Predicting Chloroform Production from Organic Precursors

Tom Bond^{1*} and Nigel Graham²

1. Department of Civil and Environmental Engineering, University of Surrey, Guildford,
GU2 7XH, UK

2. Department of Civil and Environmental Engineering, Imperial College London, London
SW7 2AZ, UK

* Corresponding author, Tel: +44(0)1483684474, email: t.bond@surrey.ac.uk

20 Abstract

21 Quantitative methods which link molecular descriptors for recognized precursors to
22 formation of drinking water disinfection byproducts are scarce. This study aimed to develop a
23 simple mathematical tool for predicting chloroform (trichloromethane) yields resulting from
24 aqueous chlorination of model organic precursors. Experimental chloroform yields from 211
25 precursors were collated from 22 literature studies from 1977 onwards. Nineteen descriptors,
26 some established and others developed during this study, were used as inputs in a multiple
27 linear regression model. The final model, calibrated using five-way leave-many-out cross-
28 validation, contains three descriptors. Two novel empirical descriptors, which quantify the
29 impact of adjacent substituents on aromatic and enolizable chlorine substitution sites, were
30 the most significant. The model has $r^2 = 0.91$ and a standard error of 8.93% mol/mol.
31 Experimental validation, using 10 previously untested precursors, showed a mean
32 discrepancy of 5.3% mol/mol between experimental and predicted chloroform yields. The
33 model gives insight to the influence that specific functional groups, including hydroxyl,
34 chlorine and carboxyl, have on chloroform formation and the relative contributions made by
35 separate substitution sites in the same molecule. It is anticipated that the detailed approach
36 can be updated and extended as new experimental data emerges, to encompass additional
37 precursors and groups of disinfection byproducts.

38 **Keywords:** trihalomethanes; THMs; disinfection byproducts; model compounds; QSAR

39 1 Introduction

40 The discovery, in the 1970s, that chloroform and other trihalomethanes are generated from
41 chlorination of natural organic matter during water treatment (Rook 1974) surprised the
42 scientific community. This breakthrough led to much research into disinfection byproducts
43 and regulations for the four chlorinated and/or brominated trihalomethanes being introduced

44 in the USA by the decade's end. The limit for total trihalomethanes is currently $80 \mu\text{g}\cdot\text{L}^{-1}$ in
45 the USA, with five haloacetic acids, another group of halogenated disinfection byproducts,
46 regulated at a total of $60 \mu\text{g}\cdot\text{L}^{-1}$. Total trihalomethanes are also regulated in the EU at 100
47 $\mu\text{g}\cdot\text{L}^{-1}$. Initial concern about the trihalomethanes was based on results of a rodent bioassay
48 which classified chloroform as a suspected human carcinogen (NCI 1976). Independently,
49 epidemiological studies have shown that long-term consumption of chlorinated drinking
50 water is associated with an enhanced risk of developing bladder cancer, although the
51 underlying reasons remain obscure (Hrudey 2009). However, analysis of more recent
52 toxicological evidence indicates that neither the trihalomethanes nor the haloacetic acids are
53 plausible bladder carcinogens at typical drinking water concentrations (Hrudey 2009). Thus,
54 other disinfection byproducts such as the halobenzoquinones (Zhao et al. 2012), various other
55 halogenated aromatic byproducts (Zhang et al. 2008), nitrosamines, haloacetonitriles and
56 haloacetamides (Shah and Mitch 2012) remain the focus of much current research attention,
57 as these may be more toxicologically-significant.

58 The trihalomethanes and haloacetic acids are essentially viewed by regulators of drinking
59 water quality as indicators of the total occurrence of chlorination disinfection byproducts
60 (USEPA 2015). They also remain the most-studied disinfection byproducts, particularly the
61 trihalomethanes. Model compounds have been heavily used to elucidate mechanistic
62 formation pathways and precursor characteristics from the early days of disinfection
63 byproduct research (Bond et al. 2012, Rook 1977). It was quickly appreciated that meta-
64 substituted aromatic compounds are reactive chloroform precursors. For example, Rook
65 (1977) reported that resorcinol (1,3-dihydroxybenzene) was converted into chloroform at a
66 85% mol/mol yield during chlorination, whereas yields from its regioisomers, catechol (1,2-
67 dihydroxybenzene) and hydroquinone (1,4-dihydroxybenzene), were far lower at 0.5% and
68 1.5% mol/mol, respectively (see Figure S2 for a simplified mechanism for chloroform

69 production from resorcinol). This variation can be (qualitatively) explained by the two
70 activating hydroxyl groups promoting electrophilic substitution reactions at the ortho- and
71 para- positions of the aromatic ring. However, the presence of additional substituents
72 complicates this pattern, as depending on their identity and position, chloroform yields can
73 either be enhanced or suppressed. Thus, de Laat et al. (1982), reported chloroform yields for
74 pyrogallol (1,2,3-trihydroxybenzene), 4-hydroxycatechol (1,2,4-trihydroxybenzene) and
75 phloroglucinol (1,3,5- trihydroxybenzene) as 0.1, 15.5 and 93 % mol/mol, respectively.
76 Certain aliphatic compounds, notably β -dicarbonyl compounds (Boyce and Hornig 1980,
77 Dickenson et al. 2008), including 3-oxopentanedioic acid (Table 1), also act as reactive
78 trihalomethane precursors. Formation of trihalomethanes from carbonyl compounds can be
79 likened to the haloform reaction, used for the synthetic preparation of trihalomethanes from
80 methyl ketones (Larson and Weber 1994). Its rate is controlled by the initial enolization of
81 the organic precursor and the mechanism proceeds via electrophilic addition of chlorine at the
82 carbon alpha to the carbonyl group (Figure S1).

83 Despite the extensive amount of research effort on this subject over the past ~40 years,
84 quantitative methods to predict disinfection byproduct formation based on molecular
85 descriptors are scarce. Hence, the aim of this study was to develop a simple mathematical
86 method for predicting chloroform yields from model compounds, with the expectation that
87 this will prove a useful screening tool for compounds which have not been tested
88 experimentally. To achieve this, descriptors are required which quantitatively link reactive
89 precursor functionalities to chloroform yields (and ultimately the key pathways leading from
90 one to the other). A secondary aim was that such a mathematical framework would illuminate
91 our knowledge of the characteristics of reactive precursors. Although quantitative structure
92 activity relationship (QSAR) and quantitative structure–property relationship (QSPR) models
93 are widely used in other fields, they have found limited use in disinfection byproduct

94 research, despite having the potential to streamline research efforts (Chen et al. 2015). While
95 not focussed on disinfection byproducts, the paper by Luilo and Cabaniss (2010) is
96 noteworthy as it details a QSPR, validated using literature on 201 organic compounds, for
97 predicting chlorine demand based on eight molecular descriptors. The same authors
98 subsequently developed a model for predicting chloroform formation from organic precursors
99 (Luilo and Cabaniss 2011b), although this used a smaller subset of 117 model compounds.

100 **2 Methods**

101 **2.1 Literature on chloroform formation**

102 Chloroform formation data from 22 studies (Bond et al. 2009, Bond et al. 2014, Bond et al.
103 2016, Boyce and Hornig 1980, Boyce and Hornig 1983, Bull et al. 2006, Chaidou et al. 1999,
104 Chang et al. 2011, de Laat et al. 1982, de Leer and Erkerlens 1988, Dickenson et al. 2008,
105 Gallard and von Gunten 2002, Hong et al. 2009, Hureiki et al. 1994, Larson and Rockwell
106 1979, Navalon et al. 2008, Norwood et al. 1980, Rook 1977, Rule et al. 2005, Tawk et al.
107 2015, Tomita et al. 1981, Westerhoff et al. 2004) spanning the years 1977 – 2016 were
108 collated and converted into units of % mol/mol where necessary. All studies measured
109 chloroform formation from organic precursors under formation potential conditions, i.e. using
110 excess chlorine. However, since there are no standard conditions for these tests, experimental
111 conditions vary in the literature (Table S1). In this study only data collected at pH 7-8 using
112 an excess of chlorine and contact times over 0.5 h were included. Thus, the modelling results
113 only apply to these conditions, which are representative of full-scale drinking water
114 chlorination. Median conditions from the studies included were pH = 7, contact time = 24 h,
115 temperature = 20 °C and chlorine dose = 20 mol/mol. One important difference with the
116 chlorination of natural waters is that model compound studies are typically undertaken in the
117 absence of bromide, and thus chloroform is the only one of the four chloro- and/or bromo-

118 trihalomethanes monitored. In contrast, ambient bromide in natural waters leads to formation
119 of varying amounts of brominated trihalomethanes.

120 For compounds tested in multiple studies, mean values were calculated and are given in
121 Table S2. For example, chloroform yields from the well-studied precursor resorcinol have
122 been reported 12 times, giving a mean value of 81.1% mol/mol and a standard deviation of
123 17.3 (Table S2). In total there were 69 compounds included with multiple chloroform yields.
124 The mean of the standard deviations for these repeated compounds is 5.1% mol/mol. The
125 final list of 211 precursors used for modelling, together with their chloroform yields,
126 structures and alternative names, is given in Table S2.

127 **2.2 Descriptor selection**

128 Three descriptors used by Luilo and Cabaniss (2011b) to model chloroform formation and
129 another three used by the same authors to model total organic halogen formation (Luilo and
130 Cabaniss 2011a) were included. Respectively these are the carbonyl index (CI), the difference
131 between the sum of strong electron-donating groups and the sum of carbonyls per carbon in
132 each molecule (EDCORH), the number of 1,3-activated aromatic carbons (OTactC)
133 (chloroform descriptors) and the number of phenolic groups per carbon (ArOH:C), the square
134 root of the number of heteroatoms (sqHeA) and the log of the hydrogen to carbon ratio (log
135 H:C) (total organic halogen descriptors). Hammett and Taft constants account for substituent
136 effects in aromatic and aliphatic compounds, respectively, and have been used widely in the
137 development of QSARs and linear free energy relationships. In this study they were used in a
138 manner following that described by Lee and von Gunten (2012) and Gallard and von Gunten
139 (2002), who showed that the sums of Hammett or Taft for organic compounds can be
140 quantitatively linked to rate constants for reactions with aqueous chlorine. Taft/Hammett
141 constants were taken from published sources (Hansch et al. 1995, Perrin et al. 1981). Four

142 descriptors were calculated by summing Taft constants for substituents around enolizable
143 functionalities (Enol Taft and Enolizable Taft), alkenes (Alkene Taft) and amino acids
144 (AmAc Taft). Hammett constants were summed to account for ortho-, meta- and para-
145 interactions in aromatic compounds, as well as total interactions (HammettOrtho, Hammett
146 Meta, Hammett Para and Hammett Sum). Finally, empirical constants were used to develop
147 novel descriptors for five important precursor categories: alkenes (Alkene Score), enolizable
148 aliphatics (Enolizable Score), aromatic ketones (Aromatic Ketone Score), beta-dicarbonyl
149 (BDicarb Score) and aromatic compounds (Aromatic Score). These were calculated by giving
150 each substituent around a potential chlorine substitution site a score, which were then
151 multiplied together to give a combined score for the substitution site. Therefore, these
152 descriptors quantify the influence of specific functional groups around a chlorine substitution
153 site on chloroform formation.

154 For more complex molecules, scores for individual substitution sites were summed to obtain
155 a total score for the whole molecule. Empirical substituent constants, used to derive a score
156 for a substitution site, were selected to minimise regression residuals, in a similar fashion to
157 some descriptors developed by Luilo and Cabaniss (2010). The 19 descriptors introduced
158 above were used as inputs in a multiple linear regression model in SPSS, with no y-intercept,
159 and with chloroform yield (% mol/mol) as the dependent variable, i.e.:

$$160 \quad CHCl_3 = \sum_{j=1}^M \beta_j x_j$$

161 Significant descriptors were then identified by successively eliminating the least significant
162 descriptor (that with the highest p value) until all remaining descriptors were significant
163 ($p < 0.05$), following previous work (Luilo and Cabaniss 2010, Luilo and Cabaniss 2011b).
164 This gave four significant descriptors: Enolizable Score, SqHeA, OTactC and Aromatic
165 Score. However, a correlation of $r = 0.833$ between Aromatic Score and OTactC indicated

166 these two descriptors were not independent (Table S5), so OTactC was removed as it had the
167 higher p value. All other pairwise correlations between these descriptors had $r < 0.7$,
168 indicating they are sufficiently independent to be used in multiple linear regression modelling
169 (Eriksson et al. 2003). The three named descriptors were used in all subsequent modelling
170 procedures and will now be described in more detail. The Enolizable Score multiplies
171 empirical constants for a carbonyl group, denoted as R1, with those for the other substituents
172 (R2 and R3) surrounding an enolizable site (Table 3). An alpha proton combined with a
173 carbonyl group is defined as enolizable, although the proton is excluded from the calculation,
174 as it is present in all enolizable precursors. This classification includes carboxylic acids, as
175 well as aldehydes and ketones. One descriptor value is calculated for each enolizable group,
176 even if multiple enolizable protons are present. Two carbonyls alpha (1,2-) or beta (1,3-) to
177 each other constitute a single enolizable group. Two carbonyls gamma (1,4-) to each other, or
178 more distant, constitute two separate enolizable groups. Assigning a single descriptor value to
179 each enolizable group is reasonable as introduction of one chlorine alpha to an enolizable
180 carbonyl promotes additional chlorine substitution at the same site (Larson and Weber 1994).
181 However, multiple enolizable groups in the same molecule are given separate scores, which
182 are then summed together. For example, for 3-acetylphenol (Table 1) the R1 group is
183 COC_6H_5 (OH meta), while R2 and R3 are both H (Table 3), so the Enolizable Score is $1 \times 1 \times$
184 $1 = 1$ (section 4.1 of the supplementary material). SqHeA is defined as the square root of the
185 number of heteroatoms present in a precursor (Luilo and Cabaniss 2011a). For 3-acetylphenol
186 there are two heteroatoms (Table 1), so SqHeA is 1.41 (Table S3). Finally, the Aromatic
187 Score descriptor gives each substituent in an aromatic precursor an empirical substituent
188 score (Table 4). These are then multiplied together to give a combined score for the whole
189 ring, which accounts for interactions between the substituents. Benzene is given a baseline
190 score of 1 and heterocycles 0.2. Depending on the substituents present and their interactions

191 these values can then either increase or decrease. The Aromatic Score for the ring in 3-
192 acetylphenol is 1 (baseline score for benzene) x 1 (OH in C1 position) x 0.7 (COCH₃ in C3
193 position) = 0.7 (section 4.3 of the supplementary material). Scores for separate aromatic
194 groups are summed to produce the final Aromatic Score value. Multiple detailed examples of
195 how these three descriptors were calculated are given in the section 4 of the supplementary
196 material.

197 **2.3 Model calibration and validation**

198 Initially, the complete dataset (n = 211) was randomly split into training data (n = 158) and
199 an external validation dataset (n = 53). Then the training data was split five times into a
200 calibration subset (n = 111) and an internal validation (cross validation) subset (n = 47), to
201 facilitate leave-many-out cross-validation. For this step, stratified data splitting was used, so
202 that each compound was used at least once for cross-validation (see the supplementary
203 material, section 7). Multiple linear regression modelling was used to obtain a separate
204 equation for each calibration subset of 111 compounds. Each of these five equations was used
205 to predict chloroform yields of the respective cross-validation subsets, with their performance
206 determined by comparing experimental and predicted values. The final predictive model was
207 obtained by averaging coefficients from these five calibration datasets. Chloroform yields of
208 the external validation dataset (n = 53) were predicted using the individual and average
209 equations. The applicability domain was assessed by calculating standardised residuals of
210 cross-validation and leverage for both training and external validation datasets. Using this
211 approach data points can be defined as outliers if they have standardised residuals above +2.5
212 or below -2.5 (Luilo and Cabaniss 2010). The warning leverage was calculated using $3k'/N$
213 where k' is the number of descriptors plus one, and N is the number of compounds used to
214 develop the model. Molecules in the training data and external validation which exceed the
215 warning leverage indicate molecules which are: (i) excessively influential in determining the

216 model parameters, and (ii), are predicted due to over-extrapolation of the model (i.e. they fall
217 outside the applicability domain).

218 **2.4 Experimental validation**

219 Experimental chloroform formation from compounds whose yields have not been reported
220 previously in literature was measured using gas chromatography with electron capture
221 detection (GC-ECD, Perkin Elmer Clarus 500 GC), a modified version of USEPA method
222 551.1(USEPA 1995) and a Restek Rxi-5 Sil MS column (30 m x 0.25 mm x 0.25 μ M). All
223 samples were prepared in duplicate at pH 7 (10 mM phosphate buffer). The method detection
224 limit for chloroform was 0.2 μ g·L⁻¹. Precursor concentrations ranged from 1 to 10 μ M (Table
225 5) and the chlorine dose was 20 mol/mol. All reagents were of at least analytical purity and
226 chlorine concentrations were measured using the DPD-FAS titration method (APHA et al.
227 2005). After 24 h contact time the residual chlorine was quenched using sodium sulphite
228 (APHA et al. 2005).

229 **3 Results**

230 **3.1 Model Performance**

231 Leave-many-out cross-validation produced five QSAR equations (Table S7), which were
232 averaged to generate the final equation shown in Table 2. The average model has $r^2 = 0.91$
233 and a standard error of 8.93% mol/mol (Table S7). Respective values obtained from
234 regression of chloroform formation across the complete dataset using the same three
235 descriptors were 0.90 and 8.87 % mol/mol (Table S7). The similarity of these values
236 indicates that the data splitting procedures and leave-many-out cross-validation did not
237 introduce any meaningful bias into the average model.

238 For comparison, chloroform formation was also modelled by the three descriptors (CI,
239 OTactC and EDCORH) used previously for this purpose by Luilo and Cabaniss (2011b),

240 although in the earlier case with a smaller dataset of 117 precursors. Using these same
241 descriptors and the 211 precursors included in the current study gave a model with $r^2 = 0.77$
242 and a standard error of 13.87% mol/mol (Table SI 11). The improved performance of the
243 current model highlights how the selected descriptors more accurately reflect how the
244 functionality around chlorine substitution sites relates to chloroform yields. The OTactC and
245 CI descriptors can only take a limited number of values (e.g. 0 or 1 for OTactC), whereas the
246 Aromatic Score and Enolizable Score descriptors can assume a range of values, depending on
247 the specific chemical identity of the substituents present (see below).

248 Mean statistical parameters for the fivefold leave-many-out cross-validation show the five
249 individual equations have high predictive power, since $r^2 > 0.6$; $q^2 > 0.5$; $0.85 < k$ or $k_0 < 1.15$
250 and $r^2 - r_0^2 / r^2 < 0.1$ (Table S8) (Golbraikh and Tropsha 2002). In this context q^2 is the predicted
251 r^2 value; with k/k_0 and r^2/r_0^2 the gradient and r^2 values with/without the y-intercept,
252 respectively. The latter two were obtained from plotting predicted versus chloroform
253 formation. The average model, when applied to the five individual cross validation datasets,
254 as well as to the training set ($n = 158$) also had high predictive power (Table S9; Figure 1).
255 Thus, it can be concluded that the predictive model is robust.

256 Chloroform yields for the 53 precursors in the external validation dataset were predicted
257 using the average model (Figure 1), as well as for the five individual equations derived during
258 cross validation (Table S10). The mean statistics from the five individual equations are also
259 indicative of high predictive power since they fulfil the criteria noted above. They also agree
260 well with equivalent statistics from the average equation (Table S10). Plotting predicted
261 chloroform yields against standardised residuals (Figure SI 6) shows that there is no pattern
262 amongst the compounds with higher residuals, again evidencing that the developed model is
263 appropriate for the dataset.

264 Based on the applicability domain analysis, compounds would be classified as outliers if: (i)
265 they have standardised residuals above +2.5 or below -2.5 and (ii) they exceed the warning
266 leverage of 0.095. Based on these criteria no precursors in either the training set or external
267 validation dataset were outliers (Figure S1). Nonetheless, there are a number of compounds
268 which fulfil one of these two conditions. In the training set there are four compounds which
269 have standardised residuals above 2.5. These are phlorizin, arg-lys-glu-val-tyr, 3-
270 oxohexanedioic acid and epigallocatechin gallate. However, in all these cases leverage values
271 are well below the warning value and thus they are not considered influential in deciding the
272 model descriptors. In the external validation dataset, tembotrione has a leverage above the
273 warning value. Thus, its predicted chloroform yield was over-extrapolated, although its
274 standardised residual still lies within the applicability domain (Figure S5). Other authors have
275 reported improved performance when predicting log transformed disinfection byproduct
276 formation using water quality parameters (Obolensky and Singer 2005). However, in this
277 study a model generated with log transformed chloroform data did not improve upon that
278 shown in Table 2 (see Table S6).

279 **3.2 Experimental model validation**

280 The 10 previously untested precursors comprised six phenols and four aliphatic carbonyls
281 (Table 5). The latter include alkene, carboxylic acid and β -dicarbonyl functionalities (Table
282 5). Predicted chloroform yields, estimated using the coefficients in Table 2, ranged from
283 0.6% to 21.3% mol/mol, whereas experimental values were from 1.6 – 37.3 % mol/mol
284 (Table 5). In general, experimental and predicted values compared well, with a mean
285 difference of 5.3% mol/mol across the 10 precursors, which is lower than the standard error
286 of 8.93% associated with use of the average model. Nonetheless, for 4-chloro-2-
287 methylphenol, experimental and predicted chloroform yields differed by 16.0% (Table 5).
288 This is a chlorinated phenol, a precursor category discussed in section 5.3.2.

289 3.3 Insights from the Model

290 3.3.1 Role of Significant Descriptors

291 The Aromatic Score and Enolizable Score descriptors are the most significant regarding the
292 model's operation ($p = 0.000$ for both) and their role is to quantify reactivity in these two
293 crucial precursor groups. While the contribution of SqHeA ($p = 0.033$) is less obvious, its
294 negative coefficient (Table 2) indicates chloroform formation from more complex precursors
295 would otherwise be typically slightly overestimated by the model. The absence of any
296 descriptors involving Taft or Hammett constants highlights their limited utility in predicting
297 trihalomethane formation. Since both have been quantitatively linked to rate constants of
298 reactions between organic pollutants and chlorine (Gallard and von Gunten 2002, Lee and
299 von Gunten 2012), in turn this indicates that the kinetics (of the initial reaction step) are not
300 strongly correlated with trihalomethane formation. Since the mechanistic routes which lead to
301 trihalomethane formation are complex and involve multiple steps this is perhaps not
302 unexpected.

303 3.3.2 Aromatic Precursors

304 Hammett constants for the substituent OH do not reflect the potency of phenols as chloroform
305 precursors, especially for meta-substituted compounds. The sigma meta value for OH is 0.12
306 and it is therefore less activating than CH_3 and NH_2 , which have sigma meta values of -0.07
307 and -0.16, respectively. Chloroform yields for resorcinol (1,3-dihydroxybenzene), 3-
308 aminophenol and meta-cresol (3-methylphenol) are 81.1 ± 17.3 %, 13.9 ± 5.9 % and 6.1 ± 0.4 %
309 mol/mol, respectively, which illustrates the difficulties associated with quantitatively linking
310 Hammett constants to trihalomethane formation. Similarly, OTactC scores for the first two of
311 these compounds are both 1. This descriptor, which applies to 1,3-activated aromatic carbons,
312 can only take a value of 0 or 1 (Luilo and Cabaniss 2011b). All precursors ($n = 28$) with a

313 1,3-activated aromatic carbon have a OTactC value of 1, even though their experimental
314 chloroform yields vary widely, from 4.6 to 98.0% mol/mol.

315 These comparisons explain why new empirical descriptors were developed during the course
316 of this study. Constants for OH, CH₃ and NH₂ in the C3 (meta) position, as used when
317 calculating the Aromatic Score descriptor, are 5.5, 0.8 and 1.3 (Table 4), which reflects the
318 fact that precursors containing a resorcinol structure are more reactive chloroform precursors
319 than the corresponding cresols or anilines. While aliphatic amine groups, e.g. in amino acids,
320 are typically protonated under water treatment conditions (e.g. the pK_a for ammonium (⁺NH₄)
321 is 9.25), when calculating Aromatic Score and Enolizable Score values no attempt is made to
322 distinguish between the protonated and non-protonated forms.

323 Another issue with using Hammett constants to explain the formation of halogenated
324 products is that in water disinfection applications they are typically used in an additive
325 manner (Gallard and von Gunten 2002, Lee and von Gunten 2012). Thus, their sum will
326 increase/decrease with an increasing number of electron-withdrawing/electron-donating
327 substituents around an aromatic ring. In turn, 4-hydroxycatechol (1,2,4-trihydroxybenzene)
328 is more activated than resorcinol, whereas phloroglucinol (1,3,5-trihydroxybenzene) is less
329 activated (in terms of the sum of their Hammett constants). However, this does not correlate
330 with trihalomethane formation (Table 1).

331 In contrast, the Aromatic Score descriptor is a product, so can either increase or decrease as
332 substituents are added to the aromatic ring: respective values for these three precursors are
333 1.80, 5.50 and 5.50 (Table 4). Another feature of the Aromatic Score descriptor is that it
334 accounts for blocked meta interactions. Aromatic precursors where a resorcinol-type structure
335 is bisected by an additional substituent tend to have low chloroform yields, as the preferred
336 chlorine substitution site is already occupied. For example, pyrogallol (1,2,3-

337 trihydroxybenzene) generates $0.8 \pm 1.4\%$ mol/mol of chloroform (Table 1). When calculating
338 the Aromatic Score descriptors blocked meta interactions, such as in pyrogallol, were not
339 considered.

340 There is one substituent which represents an exception to this rule: chlorine, as chlorine
341 substitution can either occur where the existing chlorine is located, or elsewhere in the ring.
342 Chlorinated phenols have variable chloroform yields: from $6.8 \pm 4.6\%$ for 2,4,6-
343 trichlorophenol to 98% for 4-chloro-1,3-benzenediol (Table S2). Aromatic Score values for a
344 chlorine substituent ortho (C2 or C6) or para (C4) to a strongly-activating group (OH or NH₂)
345 are higher than when more weakly-activating substituents are in the ortho- or para- positions
346 (Table 4). This indicates that additional chlorine substitution proceeds at the carbon bonded
347 to the existing chlorine. Nonetheless, there are several examples of meta-substituted
348 chlorinated phenols which are highly reactive precursors, for example, 3,5-dichlorophenol,
349 which generates 71.7% chloroform (Table S2) and chlorine has the second highest C3
350 Aromatic Score value of 2.2 (Table 4). Since both chlorine and OH are ortho/para directors
351 (although they are respectively deactivating and activating) this indicates that chlorine
352 substitution occurs elsewhere than the original chlorine group. For other precursors with
353 multiple chlorine groups, e.g. 2,3,4,6-tetrachlorophenol, it is unclear where subsequent
354 chlorination occurs. Predicted and experimental chloroform yields for 4-chloro-2-
355 methylphenol differed by 15.6% (Table 5). This indicates that the empirical descriptor values
356 used to calculate Aromatic Score values could be further optimised if additional experimental
357 chloroform yields were available.

358 Another substituent which forms an exception in the model is carboxyl (-CO₂). There are a
359 number of phenols containing a CO₂ group in addition to a resorcinol structure which
360 generate similar amounts of chloroform to resorcinol itself. For instance, chloroform yields
361 for 2,4-dihydroxybenzoic acid, 2,6-dihydroxybenzoic acid and 3,5-dihydroxybenzoic acid are

362 83.3±9.6%, 82.5±10.6% and 59.4±19.8%, respectively, versus 81.1±17.3% for resorcinol
363 (Table S2). Similarly, salicylic acid (2-hydroxybenzoic acid), 3-hydroxybenzoic acid and 4-
364 hydroxybenzoic acid all form a comparable amount of chloroform to phenol: from 3.0 to
365 6.9% for the three isomers, compared with 4.8±3.9% for phenol (Table S2). These
366 similarities are explained by decarboxylation of carboxylic acids during chlorination (Larson
367 and Rockwell 1979) and are dealt with by the model ignoring CO₂ groups in aromatic
368 precursors.

369 3.3.3 Enolizable Precursors

370 One helpful feature of the model is that, for molecules with multiple substitution sites, it
371 highlights which are primarily responsible for chloroform formation. Whereas benzaldehyde
372 is an unreactive precursor, with a chloroform yield of 0.1%, there are seven aromatic ketones,
373 all containing an acetophenone structure, including 3-acetylphenol (Table 1), with chloroform
374 yields from 10 – 45% (Table S2). Both –CHO and –COCH₃ are deactivating substituents
375 with respect to electrophilic aromatic substitution, as shown by respective Hammett meta
376 constants of 0.35 and 0.38. This indicates that chlorination occurs at the –COCH₃ group
377 rather than the aromatic ring. In turn, Aromatic Score values are also low - 0.05 for both
378 substituents in the C1 position - whereas the Enolizable Score R1 value for COC₆H₅ (no
379 substituents) is relatively high at 0.8 (Table 3). For 3-acetylphenol, the most reactive aromatic
380 ketone, the model suggests that both the aliphatic –COCH₃ group and the aromatic meta-
381 substituted moiety contribute significantly to its experimental chloroform yield of 45%. In
382 contrast, for acetophenone itself, which generates 10% mol/mol chloroform (Table 1), the –
383 COCH₃ group is the principal substitution site.

384 Cinnamic acid derivatives also contain aromatic and aliphatic functionalities, with both
385 assumed to be substitution sites. These precursors, and other aliphatic alkenes, are classified

386 as enolizable as there is a carbonyl (carboxylic acid) alpha to the alkene. Two descriptors –
387 Alkene Taft and Alkene Score – were developed to quantify chloroform production from
388 alkenes but neither are significant in the final model. This is because of the relatively small
389 number of aliphatic alkenes in the dataset and their generally modest chloroform yields, of up
390 to 23% for fumaric acid. While it is expedient for the model to treat these precursors as
391 enolizable; it is likely that chlorination actually occurs via electrophilic addition to the alkene.
392 Another simplification in the model is how amino acids are treated. Chlorination of this
393 group has been well studied (Bond et al. 2009, Hong et al. 2009, Hureiki et al. 1994) and
394 proceeds via conversion of the amino group to either a nitrile or aldehyde (Figure S4). The
395 proportion of the aldehyde and nitrile products depends on experimental conditions, but as
396 both can be chlorinated (Wyman et al. 1964), the alpha carbon is assigned a composite score
397 in the Enolizable Score descriptor. The low R1 value of 0.2 for amino acids reflects the fact
398 that chloroform yields from alpha amino acids are generally modest. The two exceptions are
399 L-tyrosine and L-tryptophan, where an activated aromatic ring can interact with the
400 enolizable carbon by resonance, as evidenced by a high R2/R3 value of 5.6 (Table 3).

401 There is only one substituent which enhances chloroform formation to a greater extent, this
402 being COCH₃, which has a R2/R3 value of 7.2 (Table 3). This only occurs in the R2/R3
403 position in certain β -dicarbonyls which are potent chloroform precursors, for example 3-
404 oxopentanedioic acid (Table 1). Molecules with three carbonyl groups around an alpha
405 proton, i.e. β -tricarboxyls, are extremely potent precursors. There are only two precursors in
406 this subset: sulcotrione and tembotrione. This explains their high Enolizable Score, 5.89 for
407 both, and chloroform formation, 91.0% and 99.0%, respectively (Table S2). Meanwhile, 3-
408 oxopentanedioic acid has two separate β -dicarbonyl groups (and enolization sites) (Table 1),
409 which explains its high Enolizable Score value of 4.68 (Table S4).

410 In contrast, R2/R3 values for other substituents are much lower, which reflects the low
411 chloroform formation of most enolizable compounds. This includes monosaccharides, which
412 the model assumes exist in the linear (aldehyde or ketone) form, rather than as the ether ring.
413 Monosaccharides with multiple enolizable groups - maltose, maltotriose and maltopentaose -
414 have higher chloroform formation, up to 18.9% for the latter. The model sums the
415 contributions from the individual enolizable groups in these precursors.

416 **4 Discussion: Future Research Directions**

417 An empirical model of the type described in this study is only as comprehensive as the
418 experimental data available. For some groups of precursors - phenols, amino acids and
419 monosaccharides - there is extensive data on how the presence of various substituents affects
420 chloroform production. The converse applies for other categories. One is the alkenes, which
421 initially react with aqueous chlorine to produce chlorohydrins (Larson and Weber 1994) and
422 eventually chloroform (Figure S3). It can be hypothesized that the presence of adjacent
423 electron-donating groups will encourage the formation of halogenated products, by making
424 the alkene more nucleophilic. Since the Enolizable (and Aromatic) Score descriptors are
425 empirical this means they can be updated as new experimental data emerges.

426 Similarly, it is anticipated the detailed approach can serve as a starting point for predicting
427 the formation of other disinfection byproducts. This includes the haloacetic acids and total
428 organic halogen, for which modified descriptors may be developed. At present there is less
429 data available for these groups than for the trihalomethanes. Nonetheless, it is clear that
430 formation of trihalomethanes is accompanied by other groups of halogenated byproducts,
431 with their relative yields dependent on the specific chemical functionality present and
432 proceeding via common precursor structures. For example, β -dicarbonyls can generate
433 significant concentrations of both trihalomethanes and haloacetic acids (Dickenson et al.

2008). Some of the high variability in chloroform yields between studies for 3-oxopentanedioic acid, $59.8 \pm 15.9\%$ mol/mol (Table 1), can be attributed to selected experimental conditions favoring one halogenated product over another. As seen from their respective R2/R3 Enolizable Score values of 3.3, 7.2 and 4.6, the substituent COO^- promotes chloroform formation less than COCH_3 or COCH_2CH_3 ; with respect to haloacetic acid formation the opposite can be postulated. Something comparable applies to L-tyrosine and L-tryptophan (Table 1), which are also known to generate significant amounts of dichloroacetic acid, trichloroacetic acid and dichloroacetonitrile (Bond et al. 2009). The model is a screening tool for organic molecules whose chloroform yields have not been experimentally tested. It is not designed for use with bulk water quality parameters (e.g. total organic carbon and ultraviolet absorbance). Nonetheless, there is scope to use analyses which provide information about the specific chemical identity of aquatic organics (e.g. gas or liquid chromatography with mass spectrometry detection) to link the model to natural waters

5 Conclusions

This study details and validates a comprehensive mathematical framework to predict the amount of chloroform produced from reactions between aqueous chlorine and organic precursors. The key findings are as follows:

- The final model, calibrated using five-way leave-many-out cross-validation, has $r^2 = 0.91$ and a standard error of 8.93% mol/mol. It contains three descriptors, the two most significant, developed specifically for this study, empirically quantify the impact of adjacent substituents on aromatic and enolizable chlorine substitution sites.
- Experimental validation, using 10 previously untested precursors, showed a mean discrepancy of 5.3% mol/mol between experimental and predicted chloroform yields.

457 • Aromatic carboxyl groups are ignored by the model, which accounts for blocked meta
458 interactions. For molecules with multiple substitution sites the model is helpful for
459 evaluating which are primarily responsible for chloroform formation. Notably, the
460 ketone side-group in acetophenone derivatives is a significant source of chloroform
461 formation.

462 **6 Acknowledgments**

463 The first author acknowledges the support of the Imperial College Junior Research
464 Fellowship scheme. Thanks also to Jineesha Mehta for helpful discussions.

465 **7 References**

466 APHA, AWWA and WEF (2005) Standard Methods for the Examination of Water and
467 Wastewater, American Public Health Association, Washington, DC.

468 Bond, T., Goslan, E.H., Parsons, S.A. and Jefferson, B. (2012) A critical review of
469 trihalomethane and haloacetic acid formation from natural organic matter surrogates.
470 Environ. Technol. Reviews 1(1), 93-113.

471 Bond, T., Henriot, O., Goslan, E.H., Parsons, S.A. and Jefferson, B. (2009) Disinfection
472 byproduct formation and fractionation behavior of natural organic matter surrogates. Environ.
473 Sci. Technol. 43(15), 5982-5989.

474 Bond, T., Mokhtar Kamal, N.H., Bonnisseau, T. and Templeton, M.R. (2014) Disinfection
475 by-product formation from the chlorination and chloramination of amines. J. Hazard. Mater.
476 278(0), 288-296.

477 Bond, T., Tang, S.C., Graham, N. and Templeton, M.R. (2016) Formation of disinfection
478 byproducts during the preparation of tea and coffee. Environ. Sci.: Water Res. Technol. 2,
479 196-205.

- 480 Boyce, S.D. and Hornig, J.F. (1980) Formation of chloroform from the chlorination of
481 diketones and polyhydroxybenzenes in dilute aqueous solution. In *Water chlorination:
482 Environmental Impact and Health Effects*; vol. 3. Jolley, R.L., Brungs, W. and Cumming, R.
483 (eds), pp. 131-140., Ann Arbor Science, Ann Arbor, MI.
- 484 Boyce, S.D. and Hornig, J.F. (1983) Reaction pathways of trihalomethane formation from the
485 halogenation of dihydroxyaromatic model compounds for humic acid. *Environ. Sci. Technol.*
486 17(4), 202-211.
- 487 Bull, R.J., Reckhow, D.A., Rotello, V., Bull, O.M. and Kim, J. (2006) Use of Toxicological
488 and Chemical Models to Prioritize DBP Research. Report 91135, Awwa Research
489 Foundation, Denver, CO.
- 490 Chaidou, C.I., Georgakilas, V.I., Stalikas, C., Saraçi, M. and Lahaniatis, E.S. (1999)
491 Formation of chloroform by aqueous chlorination of organic compounds. *Chemosphere*
492 39(4), 587-594.
- 493 Chang, H., Chen, C. and Wang, G. (2011) Identification of potential nitrogenous organic
494 precursors for C-, N-DBPs and characterization of their DBPs formation. *Water Res.* 45,
495 3753 - 3764.
- 496 Chen, B., Zhang, T., Bond, T. and Gan, Y. (2015) Development of quantitative structure
497 activity relationship (QSAR) model for disinfection byproduct (DBP) research: A review of
498 methods and resources. *J. Hazard. Mater.* 299(0), 260-279.
- 499 de Laat, J., Merlet, N. and Dore, M. (1982) Chloration de composés organiques: demande en
500 chlore et réactivité vis-a-vis de la formation des trihalométhanes. Incidence de l'azote
501 ammoniacal. Chlorination of organic compounds: Chlorine demand and reactivity in
502 relationship to the trihalomethane formation. Incidence of ammoniacal nitrogen. *Water Res.*
503 16(10), 1437-1450.

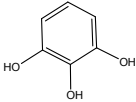
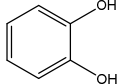
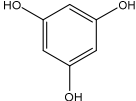
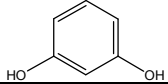
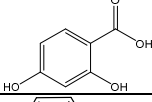
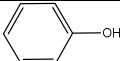
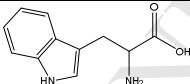
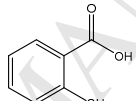
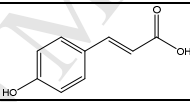
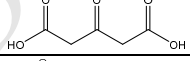
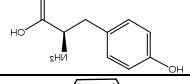
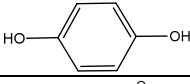
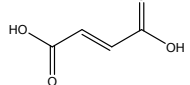
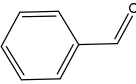
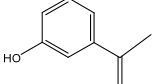
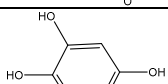
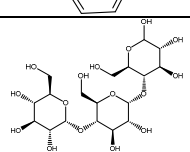
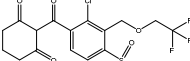
- 504 de Leer, E.W.B. and Erkerlens, C. (1988) Pathways for the production of organochlorine
505 compounds in the chlorination of humic materials. In *Biohazards of Drinking Water*
506 *Treatment*, Larson, R. A., Ed. Lewis: Chelsea, MI., 1988.
- 507 Dickenson, E.R.V., Summers, R.S., Croué, J.-P. and Gallard, H. (2008) Haloacetic acid and
508 trihalomethane formation from the chlorination and bromination of aliphatic β -dicarbonyl
509 acid model compounds. *Environ. Sci. Technol.* 42(9), 3226-3233.
- 510 Eriksson, L., Jaworska, J., Worth, A.P., Cronin, M.T.D., McDowell, R.M. and Gramatica, P.
511 (2003) Methods for reliability and uncertainty assessment and for applicability evaluations of
512 classification- and regression-based QSARs. *Environ. Health Perspect.* 111(10), 1361-1375.
- 513 Gallard, H. and von Gunten, U. (2002) Chlorination of phenols: kinetics and formation of
514 chloroform. *Environ. Sci. Technol.* 36(5), 884-890.
- 515 Golbraikh, A. and Tropsha, A. (2002) Beware of q^2 ! *J. Molecular Graphics Modell.* 20(4),
516 269-276.
- 517 Hansch, C., Leo, A. and Hoekman, D. (1995) Exploring QSAR. Hydrophobic, Electronic and
518 Steric Constants, American Chemical Society, Washington, DC.
- 519 Hong, H.C., Wong, M.H. and Liang, Y. (2009) Amino acids as precursors of trihalomethane
520 and haloacetic acid formation during chlorination. *Arch. Environ. Contam. Toxicol.* 56(4),
521 638-645.
- 522 Hrudey, S.E. (2009) Chlorination disinfection by-products, public health risk tradeoffs and
523 me. *Water Res.* 43(8), 2057-2092.
- 524 Hureiki, L., Croué, J.-P. and Legube, B. (1994) Chlorination studies of free and combined
525 amino acids. *Water Res.* 28(12), 2521-2531.
- 526 Larson, R.A. and Rockwell, A.L. (1979) Chloroform and chlorophenol production by
527 decarboxylation of natural acids during aqueous chlorination. *Environ. Sci. Technol.* 13(3),
528 325-329.

- 529 Larson, R.A. and Weber, E.J. (1994) *Reaction Mechanisms in Environmental Organic*
530 *Chemistry*, Lewis Publishers, Ann Arbor, MI.
- 531 Lee, Y. and von Gunten, U. (2012) Quantitative structure–activity relationships (QSARs) for
532 the transformation of organic micropollutants during oxidative water treatment. *Water Res.*
533 46(19), 6177-6195.
- 534 Luilo, G.B. and Cabaniss, S.E. (2010) Quantitative structure–property relationship for
535 predicting chlorine demand by organic molecules. *Environ. Sci. Technol.* 44(7), 2503-2508.
- 536 Luilo, G.B. and Cabaniss, S.E. (2011a) Predicting total organic halide formation from
537 drinking water chlorination using quantitative structure–property relationships. *SAR and*
538 *QSAR Environ. Res.* 22(7-8), 667-680.
- 539 Luilo, G.B. and Cabaniss, S.E. (2011b) QSPR for predicting chloroform formation in
540 drinking water disinfection. *SAR and QSAR Environ. Res.* 22(5-6), 489-504.
- 541 Navalon, S., Alvaro, M. and Garcia, H. (2008) Carbohydrates as trihalomethanes precursors.
542 Influence of pH and the presence of Cl⁻ and Br⁻ on trihalomethane formation potential. *Water*
543 *Res.* 42(14), 3990-4000.
- 544 NCI (1976) National Cancer Institute. Report on Carcinogenesis Bioassay of Chloroform.
545 NTIS PB-264018, Bethesda, MD.
- 546 Norwood, D.L., Johnson, J.D., Christman, R.F., Hass, J.R. and Bobenrieth, M.J. (1980)
547 Reactions of chlorine with selected aromatic models of aquatic humic material. *Environ. Sci.*
548 *Technol.* 14(2), 187-190.
- 549 Obolensky, A. and Singer, P.C. (2005) Halogen substitution patterns among disinfection
550 byproducts in the information collection rule database. *Environ. Sci. Technol.* 39(8), 2719-
551 2730.
- 552 Perrin, D.D., Dempsey, B. and Serjeant, E.P. (1981) *pKa prediction for organic acids and*
553 *bases* Chapman and Hall, New York.

- 554 Rook, J.J. (1974) Formation of haloforms during chlorination of natural water. *Water*
555 *Treatment and Examination* 23(2), 234-243.
- 556 Rook, J.J. (1977) Chlorination reactions of fulvic acids in natural waters. *Environ. Sci.*
557 *Technol.* 11(5), 478-482.
- 558 Rule, K., Ebbett, V. R. and Vikesland, P. J. (2005) Formation of chloroform and chlorinated
559 organics by free-chlorine-mediated oxidation of triclosan. *Environ. Sci. Technol.* 39(9),
560 3176-3185.
- 561 Shah, A.D. and Mitch, W.A. (2012) Halonitroalkanes, halonitriles, haloamides, and N-
562 nitrosamines: A critical review of nitrogenous disinfection byproduct formation pathways.
563 *Environ. Sci. Technol.* 46(1), 119-131.
- 564 Tawk, A., Deborde, M., Labanowski, J. and Gallard, H. (2015) Chlorination of the β -
565 triketone herbicides tembotrione and sulcotrione: Kinetic and mechanistic study,
566 transformation products identification and toxicity. *Water Res.* 76, 132-142.
- 567 Tomita, M., Manabe, H., Honma, K. and Hamada, A. (1981) Studies on trihalomethane
568 formation of model compounds by aqueous chlorination, In *Proceedings of the 8th*
569 *Symposium on Environmental Pollutants and Toxicology, vol. 28*, Sendai, Japan, 1981; pp
570 21-27.
- 571 USEPA (1995) Method 551.1. Determination of chlorination disinfection byproducts,
572 chlorinated solvents and halogenated pesticides/herbicides in drinking water by liquid-liquid
573 extraction and gas chromatography with electron capture detection. Revision 1.0, Cincinnati,
574 OH.
- 575 USEPA (2015) Stage 1 and Stage 2 Disinfectants and Disinfection Byproducts Rules. Fact
576 Sheet: Stage 2 Disinfectants and Disinfection Byproducts Rule
577 [https://www.epa.gov/dwreginfo/stage-1-and-stage-2-disinfectants-and-disinfection-](https://www.epa.gov/dwreginfo/stage-1-and-stage-2-disinfectants-and-disinfection-byproducts-rules)
578 [byproducts-rules](https://www.epa.gov/dwreginfo/stage-1-and-stage-2-disinfectants-and-disinfection-byproducts-rules) (Accessed 23 August 2016).

- 579 Westerhoff, P., Chao, P. and Mash, H. (2004) Reactivity of natural organic matter with
580 aqueous chlorine and bromine. *Water Res.* 38(6), 1502-1513.
- 581 Wyman, D.P., Kaufman, P.R. and Freeman, W.R. (1964) The chlorination of active hydrogen
582 compounds with sulfuryl chloride. II. esters, nitriles, nitro Compounds, and aldehydes. *J. Org.*
583 *Chem.* 29(9), 2706-2710.
- 584 Zhang, X., Talley, J.W., Boggess, B., Ding, G. and Birdsell, D. (2008) Fast Selective
585 Detection of Polar Brominated Disinfection Byproducts in Drinking Water Using Precursor
586 Ion Scans. *Environ. Sci. Technol.* 42(17), 6598-6603.
- 587 Zhao, Y., Anichina, J., Lu, X., Bull, R.J., Krasner, S.W., Hrudey, S.E. and Li, X.-F. (2012)
588 Occurrence and formation of chloro- and bromo-benzoquinones during drinking water
589 disinfection. *Water Res.* 46(14), 4351-4360.
- 590

Table 1: Selected chloroform precursors, taken from literature*

Name	Alternative name	Structure	Chloroform yield (% mol/mol)
Pyrogallol	1,2,3-trihydroxybenzene		0.8±1.4
Catechol	1,2-dihydroxybenzene		0.9±0.6
Phloroglucinol	1,3,5-trihydroxybenzene		74.6±22.4
Resorcinol	1,3-dihydroxybenzene		81.1±17.3
2,4-dihydroxybenzoic acid	β-resorcylic acid		83.3±9.6
Phenol			4.8±3.9
L-tryptophan			29.5±27.4
Salicylic acid	2-hydroxybenzoic acid		3.0±1.2
p-coumaric acid	4-hydroxycinnamic acid		1.3±0.6
3-oxopentanedioic acid	1,3-acetonedicarboxylic acid		59.8±15.9
L-tyrosine			16.0±16.4
Hydroquinone	1,4-dihydroxybenzene		1.7±1.1
Fumaric acid	(2E)-2-butenedioic acid		23.0
Benzaldehyde			0.1
3-acetylphenol	1-(3-hydroxyphenyl)ethanone		45.0
4-hydroxycatechol	1,2,4-trihydroxybenzene		15.5
Maltotriose			11.8
Tembotrione			99.0

*The full list, including references, is given in the supporting information (Table S2).

Table 2: Average model descriptors and standard errors obtained using leave-many-out cross-validation

	Descriptor (x_j)		
	Enolizable Score	SqHeA	Aromatic Score
Coefficient (β_j)	17.89	-0.78	14.13
Standard error	0.88	0.49	0.61

Table 3: empirical constants used to calculate the Enolizable Score descriptor

Carbonyl (R1)	Constant value	R2/R3*	Constant value
COO ⁻	0.1	H	1
CHO	0.1	OH	1.3
COCH ₃	0.6	CH ₂ OH	1.1
COCH ₂ CH ₃	0.8	OCH ₃	1.1
CONH ₂	0.05	CH ₃	0.8
COCOO ⁻	0.05	C ₆ H ₅	1
COCOCH ₂ CH ₃	0.4	OC ₆ H ₅	1
COCHOH	0.1	CH ₂ COO ⁻	1.5
CONHCH ₃	0.4	CONH ₂	1
COC ₆ H ₅ (OH para and OCH ₃ meta)	0.7	CH ₂ CH ₃	0.8
COC ₆ H ₅ (no substituents)	0.8	CH ₂ CH ₂ CH ₃	0.9
COC ₆ H ₅ (OH ortho)	0.7	CH ₂ COCH ₃	0.4
COC ₆ H ₅ (3 x OCH ₃ in 3, 4 and 5)	0.8	CH ₂ NH ₂	1.1
COC ₆ H ₅ (OH para)	1.2	NHCOCH ₃	0.8
COC ₆ H ₅ (2 x OCH ₃ meta, OH para)	1	NH ₂	1.1
COC ₆ H ₅ (OH meta)	1	Inactivated aromatic ring	1
COC ₆ H ₅ (2 x OH ortho and para, ether ortho)	1	Activated aromatic ring	5.6
COCH ₂ COO ⁻	0.6	COO ⁻	3.9
COCH ₂ COCH ₃	0.3	CH(OH)CH ₃	1.1
Amino acid	0.2	NHCH ₃	1.1
		C(OH)COO ⁻	1.2
		CH ₂ CH ₂ OH	1.1
		COCH ₂ CH ₃	4.9
		CH ₂ CHO	1.2
		CH ₂ CH ₂ COCH ₃	1.2
		COC ₆ H ₅	2 (R2) or 1.5 (R3)*
		COCH ₃	7.2
		CH(CH ₃) ₂	0.9
		CH ₂ CONH ₂	1
		SH	1.2
		CH ₂ C ₆ H ₅	1
		CH=CHNH ₂	0.8
		CH=CH ₂	1.5
		CH=(CH ₃)(COO ⁻)	1.8
		CH=CHCOCH ₃	0.8
		CH=CHCOO ⁻	1.8

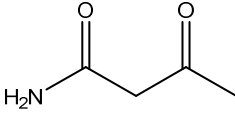
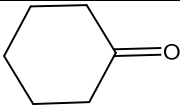
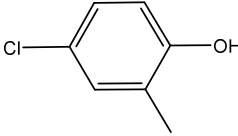
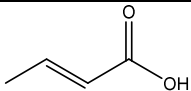
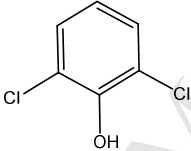
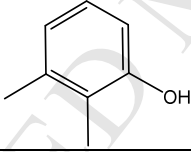
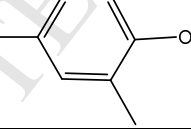
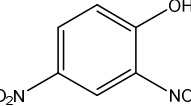
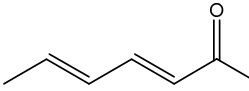
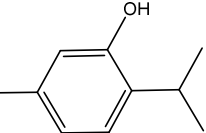
*R2 and R3 value are identical with one exception: where both are carbonyls, which only applies to sulcotrione and tembotrione.

Table 4: empirical constants used to calculate the Aromatic Score descriptor

Baseline values: 1 (benzene) or 0.1 (heterocycles)						
Substituent	C1	C2	C3	C4	C5	C6
OH	1	0.3	5.5	0.3	0.5	0.6
Cl	0.1	<i>0.6/0.9</i>	2.2	<i>0.6/1.4</i>	1.5	<i>0.5/0.7</i>
NH ₂	0.99	0.3	1.3	0.5	0.4	0.4
N(C ₂ H ₅) ₂	1.1					
NO ₂	0.2	0.2	0.8	0.2		
OCH ₃	0.3	0.5	0.5	0.5	0.5	0.2
OCH ₂ CH ₃	0.3	0.3	0.4	0.3	0.3	0.2
OCH=CH ₂	0.5	0.9	0.9	0.9	0.9	0.5
CH ₂ COO-	0.2	0.3	0.3	0.3	0.3	0.15
CH ₂ OH	0.2	0.4	0.4	0.3		0.3
CH ₃	0.3	0.85	0.8	0.8	0.8	0.5
CH ₂ CH ₃	0.3	0.5	0.7	0.5	0.7	0.4
Large alkyl group	0.3	0.4	0.6	0.4	0.5	0.3
CH=CH ₂	0.25					
COO-	Ignore (see text)					
COOCH ₃	0.05	0.05	0.05	0.05	0.05	0.05
CHO	0.05	0.2	0.4	0.05	0.05	
COCH ₃	0.05	0.6	0.7	0.6	0.6	0.6
COCH=CH ₂	0.1	0.8	0.8	0.8	0.8	0.8
OCH ₂ COO-	0.05					
CONH ₂	0.2		0.3	0.05	0.3	
CSNH ₂	0.3					
NHCOCH ₃	0.2					
fused ring	0.3	0.8	0.8	0.8	0.8	0.5
Other aryl (not fused)	0.3	0.2	0.4	0.1		
C=O (as part of ring)	0.05	0.5	0.5	0.5	0.5	0.5
CH ₂ CH(NH ₂)COO-	0.1	0.7	0.8	0.7		
CH ₂ CH ₂ CH ₂ OH	0.3					
CH=NOH	0.3					
CH=CHCHO	0.2	0.4	0.5	0.3		
CH ₂ CH ₂ COO-	0.2	0.3	0.4	0.3	0.3	
CH=CHCOO-	0.2	0.4	0.5	0.3	0.3	
CH=CHCH ₂ OH	0.05	0.4	0.5	0.3		
SO ₂ CH ₃	0.2	0.3	0.5	0.4	0.4	0.3
I				1		
CN				1		

*Chlorine (Cl) ortho (C2 or C6) or para (C4) to a strongly-activating group (OH or NH₂) takes a higher values than in the presence of other (more weakly-activating) substituents. The higher values are shown in italics above.

Table 5: Experimental chloroform yields from previously untested precursors

Name (alternative name)	Structure	Precursor concentration (μM)	Experimental (and predicted) chloroform formation (% mol/mol)
Acetoacetamide (3-oxobutanamide)		1	7.6 \pm 0.7 (9.4)
Cyclohexanone		3	2.0 \pm 0.6 (10.5)
4-chloro-2-methylphenol (4-chloro-o-cresol)		1	37.3 \pm 0.4 (21.3)
Crotonic acid (2E)-2-butenoic acid)		3	6.1 \pm 4.1 (1.6)
2,6-dichlorophenol		1	8.6 \pm 0.9 (7.6)
2,3-dimethylphenol (2,3-xylenol)		3	8.4 \pm 0.6 (13.4)
2,4-dimethylphenol (2,4-xylenol)		3	2.9 \pm 1.4 (10.9)
2,4-dinitrophenol		10	1.6 \pm 1.1 (0.6)
Sorbic acid (2E,4E)-2,4-hexadienoic acid)		1	3.8 \pm 2.2 (1.6)
Thymol (2-isopropyl-5-methylphenol)		3	10.1 \pm 0.6 (5.0)

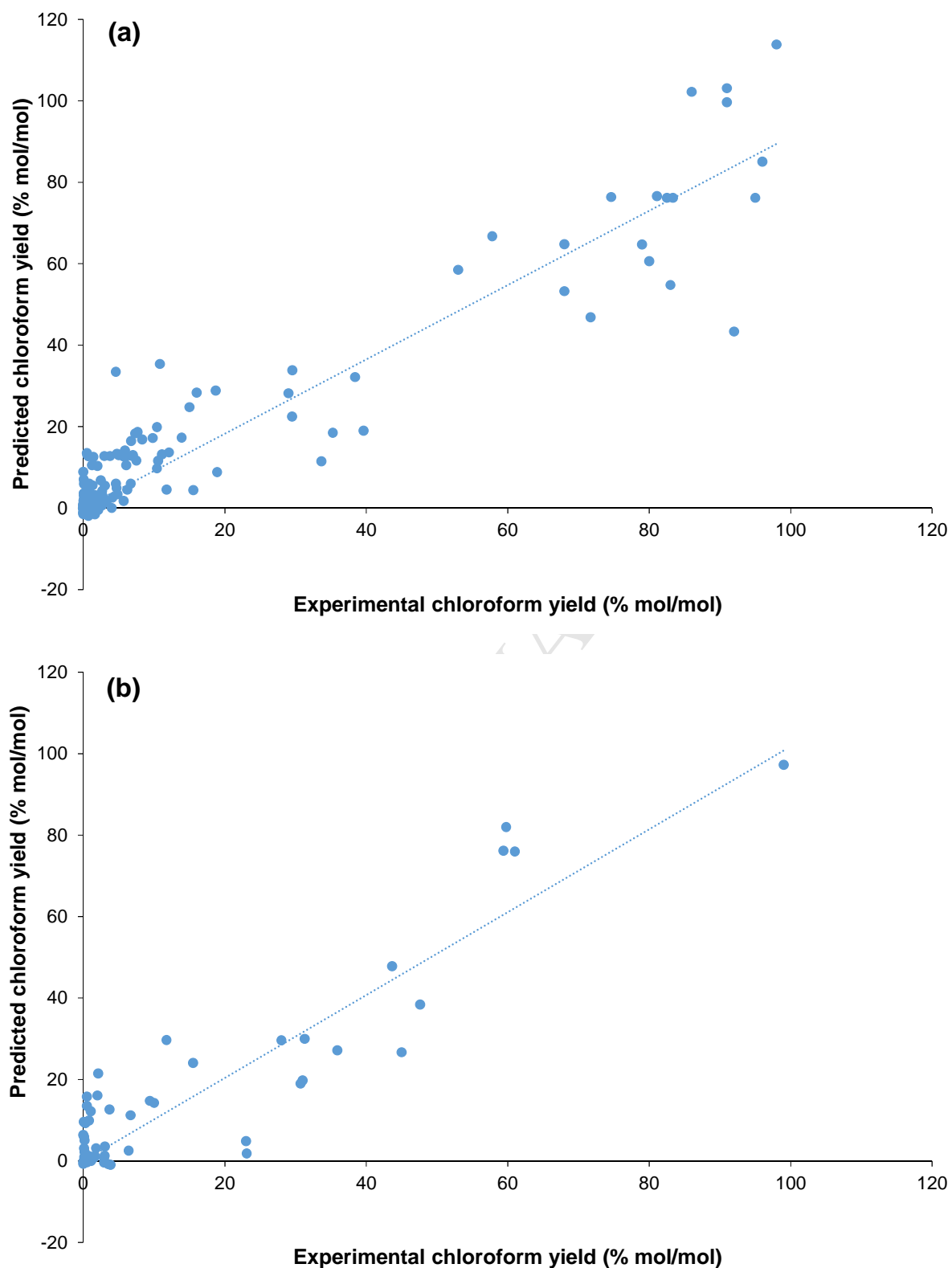


Figure 1: predicted versus experimental chloroform formation for (a) the training set (n = 158; $r^2 = 0.91$, standard error = 8.75) and (b) the external validation set (n = 53; $r^2 = 0.90$, standard error = 9.70) using the average model

Highlights

- Chloroform yields from aqueous chlorination of organic precursors were modelled.
- Novel descriptors were developed for the model.
- The final model has $r^2 = 0.91$ and a standard error of 8.93% mol/mol.
- Experimental validation was undertaken with previously unknown precursors.