

# Prediction of Retention Time and Collision Cross Section ( $CCS_{H^+}$ , $CCS_{H^-}$ , and $CCS_{Na^+}$ ) of Emerging Contaminants Using Multiple Adaptive Regression Splines

Alberto Celma,<sup>#</sup> Richard Bade,<sup>#</sup> Juan Vicente Sancho, Félix Hernandez, Melissa Humphries,\* and Lubertus Bijlsma\*



Cite This: *J. Chem. Inf. Model.* 2022, 62, 5425–5434



Read Online

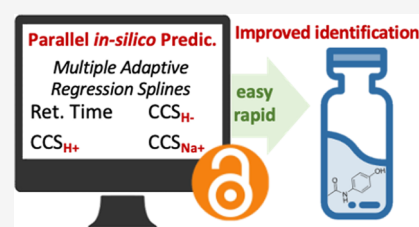
ACCESS |

Metrics & More

Article Recommendations

Supporting Information

**ABSTRACT:** Ultra-high performance liquid chromatography coupled to ion mobility separation and high-resolution mass spectrometry instruments have proven very valuable for screening of emerging contaminants in the aquatic environment. However, when applying suspect or nontarget approaches (*i.e.*, when no reference standards are available), there is no information on retention time (RT) and collision cross-section (CCS) values to facilitate identification. *In silico* prediction tools of RT and CCS can therefore be of great utility to decrease the number of candidates to investigate. In this work, Multiple Adaptive Regression Splines (MARS) were evaluated for the prediction of both RT and CCS. MARS prediction models were developed and validated using a database of 477 protonated molecules, 169 deprotonated molecules, and 249 sodium adducts. Multivariate and univariate models were evaluated showing a better fit for univariate models to the experimental data. The RT model ( $R^2 = 0.855$ ) showed a deviation between predicted and experimental data of  $\pm 2.32$  min (95% confidence intervals). The deviation observed for CCS data of protonated molecules using the  $CCS_H$  model ( $R^2 = 0.966$ ) was  $\pm 4.05\%$  with 95% confidence intervals. The  $CCS_H$  model was also tested for the prediction of deprotonated molecules, resulting in deviations below  $\pm 5.86\%$  for the 95% of the cases. Finally, a third model was developed for sodium adducts ( $CCS_{Na^+}$ ,  $R^2 = 0.954$ ) with deviation below  $\pm 5.25\%$  for 95% of the cases. The developed models have been incorporated in an open-access and user-friendly online platform which represents a great advantage for third-party research laboratories for predicting both RT and CCS data.



## 1. INTRODUCTION

In the last decade, considerable effort has been devoted to enhance the performance of high-resolution mass spectrometry (HRMS) suspect screening (SS) and nontarget screening (NTS) strategies.<sup>1–3</sup> The instrumental improvements of HRMS instruments have required the development of more sophisticated algorithms to be able to handle the large amount of data generated.<sup>3,4</sup> Therefore, the development of open-access scripts for data processing and *in silico* prediction tools represents a step-forward into the applicability of SS and NTS in wide-scope campaigns by facilitating the identification process.<sup>5–7</sup> Furthermore, the establishment of community-adopted levels of confidence for the identification of compounds using chromatography coupled to HRMS has been of paramount importance for the comparison of data across studies.<sup>8</sup>

Recently, ion mobility separation (IMS) coupled to HRMS instruments (IMS-HRMS) has proven promising for SS and NTS strategies.<sup>9</sup> It permits, in theory, to resolve co-eluting compounds with the same nominal or exact mass that could not be previously separated with solely the chromatographic method, such as isobaric or isomeric compounds.<sup>9–11</sup> Moreover, it allows the removal of mass spectrometric peaks that do not correspond to the feature of interest, which is particularly beneficial in data independent acquisition (DIA) experi-

ments.<sup>9,10,12</sup> As a consequence, there is a reduction in the necessity of data-dependent analysis because the full-spectrum HRMS acquisition can be filtered on both RT and ion mobility data.<sup>12,13</sup>

Collision cross section (CCS) values, derived from drift time (DT) measured by IMS, are known to be system- and matrix-independent and, therefore, experimental CCS data can be included in home-made or online databases with an expected deviation below 2% for most cases.<sup>9,14,15</sup> However, this is not the case for absolute chromatographic retention times (RT) which cannot easily be compared between instrumental configurations even when RT indexing approaches are applied.<sup>16–18</sup> Thus, reference standards are practically required for building a home-made database. However, SS and NTS strategies for the identification of emerging contaminants are commonly applied prior to the acquisition of the corresponding reference standards

Received: July 6, 2022

Published: October 24, 2022



<sup>1,3</sup> and, therefore, lacking any information on experimental RT and CCS. In this sense, in silico prediction tools of either chromatographic retention data or ion mobility data are of great utility to decrease the number of candidates to investigate and, therefore, increase the chance of correct identification of features.<sup>6</sup>

Several studies have predicted RT,<sup>6,19–28</sup> CCS values,<sup>7,29–36</sup> or both.<sup>13</sup> Predictors of RT have been developed mainly to model RT data in reverse-phase liquid chromatography (RPLC) and hydrophilic interaction liquid chromatography (HILIC) with prediction accuracy between approximately  $\pm 1$  to  $\pm 3$  min (up to 22% of the chromatographic gradient length). However, there is no clear agreement in the literature on how to express the prediction accuracy of the models or which should be the most appropriate statistical descriptor representing the prediction power of the system developed.<sup>6</sup> Although CCS could be theoretically modeled from the three-dimensional and chemical structure using supercomputing systems,<sup>34,37–39</sup> data-driven predictive models have also been developed showing predictive accuracies in the range of 3–6% for Traveling Wave Ion Mobility instruments (TWIMS)<sup>29,31,32,35</sup> and Drift Tube Ion Mobility instruments (DTIMS).<sup>30,31,33</sup> Similar prediction accuracy was obtained by Mollerup et al. in their study for the simultaneous prediction of RT and CCS.<sup>13</sup> However, these data-driven models were fed with data generated using different instruments depending on the output parameter. For the RT prediction, they used data gathered from an ultra-high performance liquid chromatography (UHPLC)-HRMS instrument, while for CCS prediction, they modeled CCS data generated with a UHPLC-IMS-HRMS instrument. Because RT variations could probably be observed across instruments, the utility of predicted RT in the identification of UHPLC-IMS-HRMS features is limited.

In general, the reported models were based on univariate or multivariate regressions,<sup>24,35</sup> artificial neural networks (ANNs),<sup>13,22,25,29,31</sup> quantitative structure-retention relationships (QSRR),<sup>6,21,40</sup> supported vector regression (SVR),<sup>30,33,36</sup> or statistical analysis.<sup>32,35</sup> Although Multivariate Adaptive Regression Splines (MARS) have been previously explored for RT prediction, there has been no prior exploration of the simultaneous prediction of RT and CCS.<sup>41,42</sup> MARS is a multivariate nonparametric regression procedure that was first proposed by Friedman.<sup>43</sup> One of the biggest advantages of MARS compared to the “black box” methods of ANNs is that they yield a straightforward model with simple quadratic relationship and, therefore, they are easy to interpret, with the interactions between variables clearly indicated.<sup>41</sup> Additionally, the developed MARS models for predicting analytically relevant parameters requires limited informatics resources and knowledge of prediction software tools and can consequently easily be performed. In this sense, MARS has previously been applied in the chemical sciences for quantitative structure-retention relationships.<sup>41,44</sup> However, the application of MARS for the combined prediction of chromatographic and ion mobility data of emerging contaminants has not previously been evaluated and reported in the literature.

In this work, a prediction model for both RT and CCS has been developed using MARS for the identification of candidates in SS and NTS strategies using UHPLC-IMS-HRMS. To facilitate other laboratories implementing this predictive tool in their workflows, a free online-available application has been released. This is, to best of the authors knowledge, the first application of MARS for the prediction of RT and CCS data.

Additionally, it is the first parallel RT and CCS predictive model for the same instrument facilitating the identification process of emerging contaminants in SS and NTS strategies.

## 2. MATERIALS AND METHODS

**2.1. Chemicals and Materials.** A set of 556 reference standards encompassing illicit drugs, hormones, mycotoxins, new psychoactive substances, pesticides, and pharmaceuticals was injected for the development of a CCS and RT library.<sup>9</sup> Table S1 of the Supporting Information shows the complete set of compounds used in the study with their SMILES (simplified molecular-input line-entry system) representation and measured RT and CCS data. This database is also available on the Zenodo online repository.<sup>45</sup> Within this data set, 477 protonated adducts ( $[M + H]^+$ ), 169 deprotonated adducts ( $[M - H]^-$ ), and 249 sodium adducts ( $[M + Na]^+$ ) were used for the development and validation of the CCS predictive models.

**2.2. Instrumentation.** Retention time and CCS data were obtained with a Waters Acquity I-Class UPLC system (Waters, Milford, MA, USA) coupled to a VION IMS-QTOF mass spectrometer (Waters, Milford, MA, USA), using an electrospray ionization (ESI) interface operating in positive and negative ionization mode and following the method presented in Celma et al.<sup>9</sup>

The chromatographic column used was a CORTECS C18 2.1  $\times$  100 mm, 2.7  $\mu$ m fused core column (Waters) at a flow rate of 300  $\mu$ L  $\text{min}^{-1}$ . Gradient elution was performed using H<sub>2</sub>O (A) and MeOH (B) as mobile phases, both with 0.01% formic acid. The percentage of B was initially set to 10%, and it was immediately linearly increased to 90% over 14 min, followed by a 2 min isocratic period, and then returned to initial conditions (at 16.1 min) with a 2 min equilibration of the column. The total run time was 18 min. The injection volume was 5  $\mu$ L.

A capillary voltage of 0.8 kV and cone voltage of 40 V were used. The desolvation temperature was set to 550  $^{\circ}\text{C}$ , and the source temperature to 120  $^{\circ}\text{C}$ . Nitrogen was used as drying and nebulizing gas. The cone gas flow was 250 L  $\text{h}^{-1}$  and desolvation gas flow of 1000 L  $\text{h}^{-1}$ . The column temperature was set to 40  $^{\circ}\text{C}$  and the sample temperature to 10  $^{\circ}\text{C}$ . MS data were acquired using the VION in HDMSe mode, over the range  $m/z$  50–1000, with N<sub>2</sub> as the drift gas, an IMS wave velocity of 250  $\text{m s}^{-1}$ , and wave height ramp of 20–50 V. Leucine enkephalin ( $m/z$  556.27658 and  $m/z$  554.26202) was used for mass correction in positive and negative ionization modes, respectively. Two independent scans with different collision energies were acquired during the run: a collision energy of 6 eV for low energy (LE) and a ramp of 28–56 eV for high energy (HE). A scan time of 0.3 s was set in both LE and HE functions. Nitrogen ( $\geq 99.999\%$ ) was used as collision-induced dissociation (CID) gas. All data were examined using an in-house built accurate mass screening workflow within the UNIFI software (version 1.9.4) from Waters Corporation.

**2.3. Retention Time and Collision Cross-Section Modeling.** **2.3.1. Molecular Descriptors.** A total of 1666 molecular descriptors were downloaded from Dragon v5.4 integrated within OChem website (Online Chemical Database with modeling environment, [www.ochem.eu](http://www.ochem.eu)).<sup>46</sup> The complete set of descriptors for the molecules used in the study is available in Table S1.

**2.3.2. Prediction Model.** MARS analysis was applied to predict both RT and CCS for protonated adducts ( $[M + H]^+$ ) in a single multivariate model. Additionally, univariate models for individual RT and CCS for protonated adducts ( $[M + H]^+$ )

( $\text{CCS}_{\text{H}}$ ) and sodium adducts ( $[\text{M} + \text{Na}]^+$ ) ( $\text{CCS}_{\text{Na}}$ ) were also performed. Because of the expected low correlation between RT and CCS ( $r = 0.354$ ), a multivariate model was not considered essential. As a further justification for this decision, the cross-validated  $R^2$  values for the multivariate model were 0.798 for RT and 0.964 for  $\text{CCS}_{\text{H}}$ . This suggests instability on the data that is varying the accuracy of the model fits (particularly for RT). Therefore, the development of a multivariate MARS model able to predict simultaneously RT and CCS simultaneously was discarded.

MARS was able to select the most suitable molecular descriptors for each model (Table 1), and predictive interval

**Table 1. Descriptors Needed for Each of the Univariate MARS Models for RT and  $\text{CCS}_{\text{H}}$  and  $\text{CCS}_{\text{Na}}$ <sup>a,b</sup>**

molecular descriptors		
RT	$\text{CCS}_{\text{H}}$	$\text{CCS}_{\text{Na}}$
ALOGP	AMR	AMR
ALOGPS_LogP	L1m	Har2
BEHm4	LPRS	MAXDN
GATS1m	MDDD	Mor17m
Mor16m	nRCHO	nR09
N-068	PCR	piID
nDB	Whetp	QXXv
nRNHR		QZZm
O-059		RDF065v
STN		ZM1v

<sup>a</sup>Note that there are no similarities between the three univariate models. <sup>b</sup>ALOGP: Ghose-Crippen octanol–water partition coefficient (logP) (calculation based on Viswanadhan et al.;<sup>49</sup>ALOGPS\_LogP: Ghose-Crippen octanol–water partition coefficient (logP) (calculation based on Tetko and Tanchuk;<sup>50</sup>AMR: Ghose-Crippen molar refractivity; BEHm4: highest eigenvalue n. 4 of Burden matrix/weighted by atomic masses; GATS1m: Geary autocorrelation – lag 1/weighted by atomic masses; Har2: square reciprocal distance sum index; L1m: 1st component size directional WHIM index/weighted by atomic masses; LPRS: log of product of row sums; MAXDN: maximal electrotopological negative variation; MDDD: mean distance degree deviation; Mor16m: 3D-MoRSE – signal 16/weighted by atomic masses; Mor17m: 3D-MoRSE – signal 17/weighted by atomic masses; N-068: Al3-N atom-centered fragment; nDB: number of double bonds; nR09: number of 9-membered rings; nRCHO: number of (aliphatic) aldehydes; nRNHR: number of secondary (aliphatic) amines; O-059: Al-O-Al atom-centered fragment; PCR: ratio of multiple path count over path count; piID: conventional bond-order ID number; QXXv: Qxx COMMA2 value/weighted by atomic van der Waals volumes; QZZm: Qzz COMMA2 value/weighted by atomic masses; RDF065v: radial distribution function – 6.5/weighted by atomic van der Waals volumes; STN: spanning tree number (log); Whetp: Wiener-type index from polarizability weighted distance matrix; ZM1v: first Zagreb index by valence vertex degrees.<sup>51</sup>

bands were constructed for the univariate cases assuming a linear model variance structure. To meet this assumption, the square root of RT was modeled. The selection of molecular descriptors was automatically performed by the MARS algorithm during the model development, so no chemical bias from the analysts would influence the method.

The  $\text{CCS}_{\text{H}}$  prediction model was also explored for the prediction of CCS for deprotonated adducts ( $[\text{M} - \text{H}]^-$ ) and sodium adducts ( $[\text{M} + \text{Na}]^+$ ).  $\text{CCS}_{\text{H}}$  accurately modeled  $[\text{M} - \text{H}]^-$  data but could not predict data at acceptable levels of accuracy for  $[\text{M} + \text{Na}]^+$ . Therefore, an exclusive univariate

model was considered for the prediction of CCS data for sodium adducts ( $\text{CCS}_{\text{Na}}$ ).

All analyses were complete using R,<sup>47</sup> and MARS analysis was completed using the earth package with the variance structure defined using the linear model (lm) option.<sup>48</sup>

### 3. RESULTS AND DISCUSSION

#### 3.1. Development and Validation of Prediction Models. 3.1.1. Individual RT and CCS Model Development.

There is no assumption of an underlying variance structure with the multivariate MARS analysis, and there was no facility to define one within the earth package at the time of implementation. However, for the univariate analyses, a linear model variance structure was defined. This meant the standard deviation was estimated as a function of the predicted response and, hence, allowed for the construction of prediction intervals.

It is essential to use prediction intervals, rather than confidence intervals, in cases where the goal is to predict future values. A prediction interval is wider than a confidence interval and, at the 95% level, will provide bounds within which 95% of predicted values should fall.

All analyses considered the whole set of 1666 molecular descriptors as possible inputs to be used in the models. The assumptions of normality, linearity, and homoscedasticity were assessed for the univariate models which held those assumptions. The univariate MARS fit to RT violated the assumptions of linearity and homoscedasticity, so a square root transform was applied. This then reasonably met assumptions.

In summary, three different univariate models were developed for the prediction of RT (eq 1), CCS data for (de)protonated molecules ( $\text{CCS}_{\text{H}}$ ) (eq 2), and CCS data for sodium adducts ( $\text{CCS}_{\text{Na}}$ ) (eq 3). As an example and to assist with interpretation, in eq 1, the term  $0.099 \cdot \max(0, (\text{nDB} - 3))$  is equal to 0 for  $\text{nDB} \leq 3$  and equal to  $0.099 \cdot (\text{nDB} - 3)$  for  $\text{nDB} > 3$ .

RT model (eq 1):

$$\begin{aligned} \sqrt{\text{RT}} = & 2.343 - 0.171 \cdot \max(0, (4.22 - \text{ALOGPS\_logP})) \\ & + 0.099 \cdot \max(0, (\text{nDB} - 3)) - 0.086 \cdot \\ & \max(0, (3 - \text{nDB})) - 0.451 \cdot \\ & \max(0, (\text{N.068} - 1)) + 0.725 \cdot \\ & \max(0, (1 - \text{N.068})) + 0.632 \cdot \\ & \max(0, (1 - \text{nRNHR})) - 2.177 \cdot \\ & \max(0, (\text{BEHm4} - 3.582)) - 0.533 \cdot \\ & \max(0, (3.582 - \text{BEHm4})) - 1.565 \cdot \\ & \max(0, (\text{Mor16m} - 0.54)) + 0.111 \cdot \\ & \max(0, (\text{ALOGP} - 2.719)) - 0.234 \cdot \\ & \max(0, (2.719 - \text{ALOGP})) + 0.114 \cdot \\ & \max(0, (\text{O.059} - 1)) - 0.138 \cdot \\ & \max(0, (1 - \text{O.059})) - 3.185 \cdot \\ & \max(0, (\text{GATS1m} - 1.422)) - 0.132 \cdot \max(0, \\ & (\text{STN} - 6.985)) \end{aligned} \quad (1)$$

CCS<sub>H</sub> model (eq 2):

$$\begin{aligned} \text{CCS}_{\text{H}} = & 203.344 + 0.482 \cdot \max(0, (\text{AMR} - 94.347)) \\ & - 0.524 \cdot \max(0, (94.347 - \text{AMR})) \\ & - 0.002 \cdot \max(0, (\text{Whetp} - 1940.49)) \\ & - 0.836 \cdot \max(0, (9.95 - \text{L1m})) - 14.618 \\ & \cdot \max(0, (\text{PCR} - 1.109)) + 36.31 \\ & \cdot \max(0, \text{nRCHO}) + 0.361 \cdot \\ & \max(0, (\text{LPRS} - 171.967)) - 0.157 \\ & \cdot \max(0, (171.967 - \text{LPRS})) - 0.74 \\ & \cdot \max(0, (28.622 - \text{MDDD})) \end{aligned} \quad (2)$$

CCS<sub>Na</sub> model (eq 3):

$$\begin{aligned} \text{CCS}_{\text{Na}} = & 197.356 - 0.252 \cdot \max(0, (102.616 - \text{AMR})) \\ & + 0.575 \cdot \max(0, (\text{Har2} - 117.656)) \\ & - 0.793 \cdot \max(0, (117.656 - \text{Har2})) \\ & - 5.873 \cdot \max(0, (\text{nR09} - 1)) - 5.475 \\ & \cdot \max(0, (1 - \text{nR09})) + 0.046 \\ & \cdot \max(0, (158.403 - \text{QXXv})) + 0.074 \\ & \cdot \max(0, (527.605 - \text{ZM1V})) - 0.038 \\ & \cdot \max(0, (470.721 - \text{QZZm})) + 8.192 \\ & \cdot \max(0, (\text{Mor17m} - 0.302)) + 12.649 \\ & \cdot \max(0, (-0.302 - \text{Mor17m})) + 0.116 \\ & \cdot \max(0, (171.057 - \text{piID})) + 1.392 \\ & \cdot \max(0, (\text{MAXDN} - 2.491)) + 4.682 \\ & \cdot \max(0, (2.491 - \text{MAXDN})) - 0.442 \\ & \cdot \max(0, (\text{RDF065v} - 6.402)) \end{aligned} \quad (3)$$

The univariate models obtained a cross-validated  $R^2 = 0.855$  for the RT model,  $R^2 = 0.966$  for the CCS<sub>H</sub> model, and  $R^2 = 0.954$  for the CCS<sub>Na</sub> model. Table 1 reveals that the univariate models (RT and CCS<sub>H</sub>) do not share a single descriptor, lending weight toward the argument that univariate models provide better fits to the data than the previously explored multivariate model.

**3.1.2. Cross-Validation of the RT, CCS<sub>H</sub>, and CCS<sub>Na</sub> Models.** MARS models were fitted using a threefold cross-validation with 30 iterations. This procedure splits the data into three sections, fits the model to two of those sections (*training data*), and then tests the accuracy of the resulting model on the final section (*test data*). This procedure is then repeated 30 times, each time randomly dividing the data into three sections. The measure of accuracy used to assess goodness of fit is the cross-validated  $R^2$ , which looks at the average  $R^2$  value obtained across all 30 iterations when the model was fit to the test data. This value is usually lower than the  $R^2$  for the best model fit but dramatic changes suggest volatility in the data or overfitting in the modeling procedure.

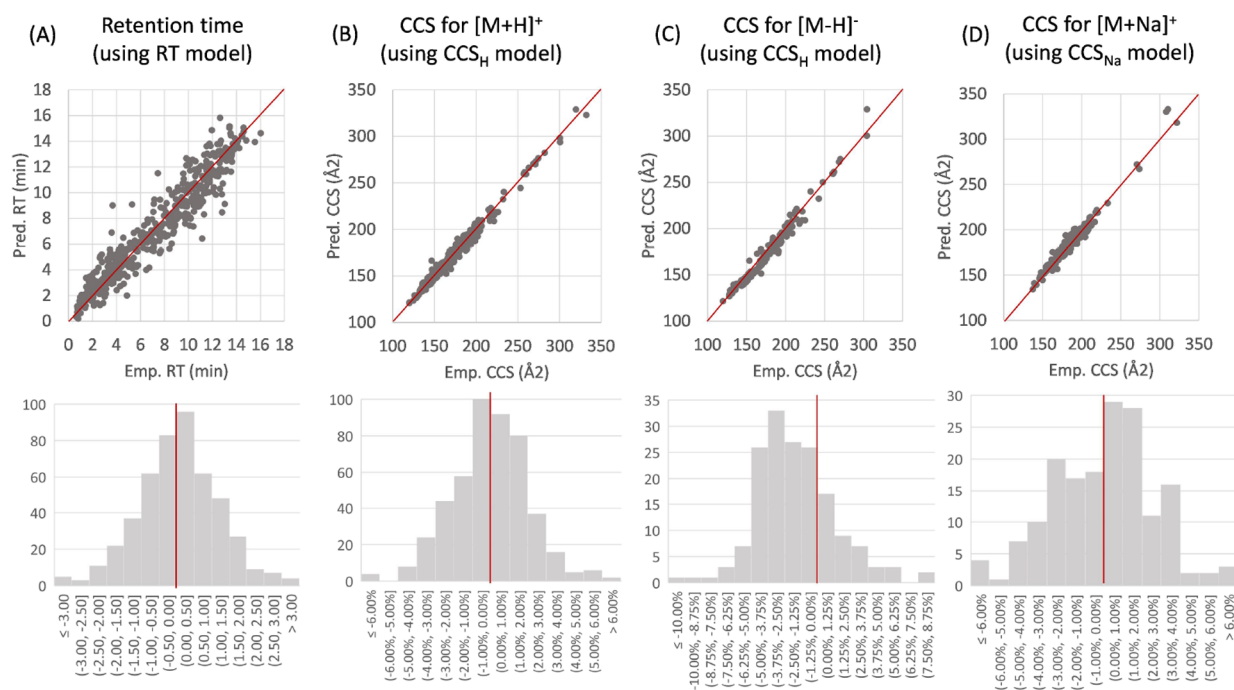
In order to perform an additional model evaluation and to obtain an overview of the model performance, RT and CCS data

were predicted for the molecules used for model development, but flagged as “unknown” substances. By comparing predicted and experimental RT data (Figure 1A, top), it was observed that the average deviation obtained using the RT model (eq 1) was  $\pm 0.72$  min, as shown in Table 2. However, 95% of the predictions fell within  $\pm 2.32$  min. Additionally, it could also be observed that deviations in predicted data distributed normally around 0% deviation (marked as a red line in Figure 1A, bottom). The prediction accuracy obtained is an improvement for the 95% intervals in previously developed models ( $\pm 4.0$  min (22%) using linear correlation *logKow* predictor,<sup>24</sup>  $\pm 2.80$  min (15%) using ANNs<sup>25</sup> over a total chromatographic run of 18 min) and in line with the model developed by means of ANN by Mollerup et al. (over  $\pm 2$  min (13%) deviation in a total run of 15 min).<sup>13</sup> The developed model herein presented also improves the prediction accuracy compared to Barron and McEneff where they obtained an average deviation of  $\pm 1.02$  min<sup>22</sup> (3–13% for the gradient length ranging 8–35 min). As another way of presenting prediction accuracy, Figure 2 plots the predicted vs experimental data with the 95% prediction intervals (blue area) for the univariate MARS analysis of the  $\sqrt{\text{RT}}$ . Approximately, only 8% of predicted RT were more than 2 min away from experimental ones. In this figure, we can also observe the 95% interval boundaries of the predicted values. This should be estimated depending on the RT because the prediction intervals are not constant across the whole chromatogram.

Prediction accuracy for CCS data was also studied. The deviation observed for CCS data of  $[\text{M} + \text{H}]^+$  using CCS<sub>H</sub> model averaged  $\pm 1.23\%$ , being  $\pm 4.05\%$  within 95% of the cases (Table 2). Figure 1B, bottom shows that deviations randomly distributed around 0% (marked as a red line) value without biasing predicted data. When compared with previous models, the CCS<sub>H</sub> model outstands the performance of developed ANNs prediction tools for CCS data of protonated molecules, which showed an accuracy of  $\pm 5$ –6% for 95% of the cases<sup>13,29</sup> or roughly  $\pm 2.5\%$  deviation for 50% of the cases.<sup>31</sup> This vast improvement in the accuracy could be explained by the larger database used for the model development as well as the better fitting of experimental data with MARS than ANNs. In addition, the present method also improves other machine learning models, such as CCSbase, which yield an accuracy slightly over  $\pm 5\%$  deviation (95% confidence interval).<sup>32</sup> The recently developed model AllCCS used more than 5000 experimental CCS values to train a support vector regression-based prediction model, which resulted in an accuracy of  $\pm 4\%$  for 84% of the cases.<sup>36</sup> The obtained accuracy is in line with that obtained in the present study, although the CCS<sub>H</sub> model is slightly more accurate for predictions because 95% of the cases have a deviation of  $\pm 4.05\%$ .

Figure 3A shows the 95% prediction intervals (blue area) for the univariate MARS analysis on the CCS<sub>H</sub> model. The blue lines are placed at predicted values  $\pm 2 \text{ \AA}^2$  and the purple are  $\pm 5 \text{ \AA}^2$ . It is clear that the model still predicts well at higher values because all data points are below the purple lines. However, because there are less data at higher CCS values, the prediction intervals are much larger to accommodate the uncertainty than they are in the low CCS values range where there are more data, resulting in a better fit.

Additionally, the application of the CCS<sub>H</sub> model for the prediction of CCS values for deprotonated molecules was tested, yielding highly accurate predictions (Figure 1C, top). By predicting mobility data for a set of 169 molecules ionized in negative mode, it was observed that the differences between the



**Figure 1.** Top: Comparison of experimental and predicted RT data using the RT model (A), CCS for protonated molecules using  $\text{CCS}_H$  model (B), CCS for deprotonated molecules using the  $\text{CCS}_H$  model (C), and CCS for sodium adducts using the  $\text{CCS}_{Na}$  model (D). (Red line indicates region where Experimental CCS = Predicted CCS) Bottom: Histogram distribution of deviations between experimental and predicted data for RT data using the RT model (A), CCS for protonated molecules using the  $\text{CCS}_H$  model (B), CCS for deprotonated molecules using the  $\text{CCS}_H$  model (C), and CCS for sodium adducts using the  $\text{CCS}_{Na}$  model (D). (Red vertical lines indicate 0% deviation).

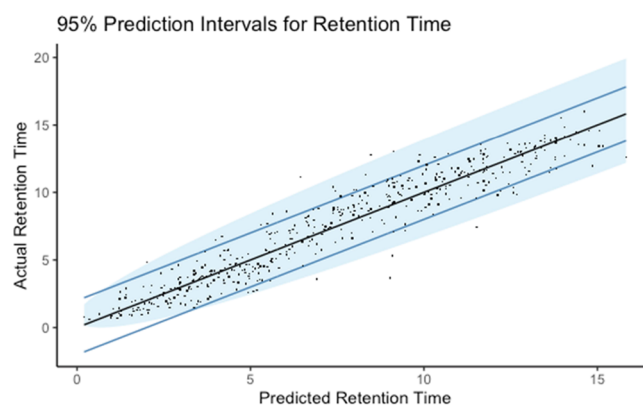
observed and predicted CCS for the  $[M - H]^-$  fell, 95% of the time, within  $-13.4$  and  $9.3 \text{ \AA}^2$ , with a slight tendency to under-predict CCS values (Figure 1C, bottom). In relative terms, average deviation for  $[M - H]^-$  data was  $\pm 2.79\%$  ( $\pm 5.86\%$  for the 95% of the cases, Table 2). Although these deviations seem larger than the ones observed for  $[M + H]^+$  data, this increase in the deviations observed for  $[M - H]^-$  was expected because the model was developed with  $[M + H]^+$  data. However, it was assumed that the predictions of the  $\text{CCS}_H$  model developed with  $[M + H]^+$  data could also be extrapolated to the prediction of CCS data for  $[M - H]^-$ , as no remarkable improvement was expected if a model was exclusively developed for deprotonated molecules.

Ideally, a unique model for the prediction of CCS for (de)protonated molecules and sodium adducts was intended.

**Table 2.** Deviations at Percentiles 50 (Average), 95, and 99 for the Predicted RT and CCS Data during Model Validation

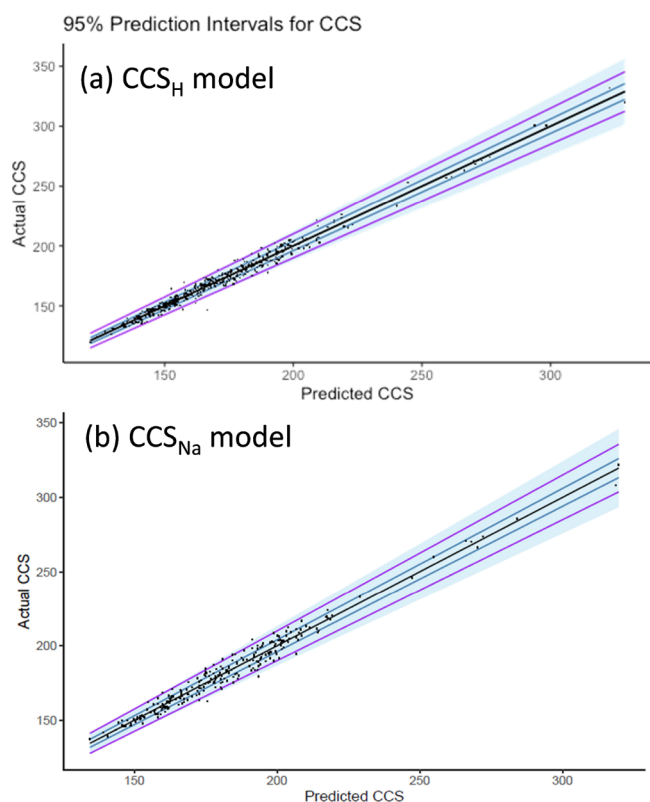
model	average deviation	deviation at 95%	deviation at 99%
RT	$\pm 0.72$ min	$\pm 2.32$ min	$\pm 3.82$ min
$\text{CCS}_H$	$[M + H]^+$	$\pm 1.23\%$	$\pm 4.05\%$
	$[M - H]^-$	$\pm 2.79\%$	$\pm 5.86\%$
	$[M + Na]^+$	$\pm 4.77\%$	$\pm 10.86\%$
$\text{CCS}_{Na}$	$[M + Na]^+$	$\pm 2.08\%$	$\pm 5.25\%$

Therefore, the  $\text{CCS}_H$  model was also tested against  $[M + Na]^+$  data. However, high deviations were observed ( $\pm 4.77\%$  average,  $\pm 10.86\%$  for the 95% of the cases, Table 2), which could be expected because of the likely higher impact of the volume of the sodium atom in the overall CCS of the molecule. In light of these data,  $[M + Na]^+$  data required a separate model for CCS prediction that was different to the one initially developed. The



**Figure 2.** 95% prediction intervals (blue area) for the univariate MARS analysis on the square root of RT. Blue lines are placed at the predicted values  $\pm 2$  min. Approximately, only 8% of observed retention times were more than 2 min away from their predicted value.

procedure for the  $\text{CCS}_{Na}$  model development was equivalent to the process described above (section 2.3) but using as input a data set of 249 CCS values for  $[M + Na]^+$  ions. The accuracy of the model was evaluated by also comparing predicted and experimental data (Table 2). Prediction deviations were  $\pm 2.08\%$  on average ( $\pm 5.25\%$  for the 95% of the cases), showing a great improvement compared with predicted data using the  $\text{CCS}_H$  model. Figure 3B depicts the predicted vs experimental CCS values comparing the 95% prediction intervals (blue colored area) for the univariate MARS analysis on the  $\text{CCS}_{Na}$  model. The fact that different predicted values can be obtained for both protonated molecules and sodium adducts is of great help for experimental observations of both species for a suspect substance. Hence, increased confidence on the tentative



**Figure 3.** 95% prediction intervals (blue area) for the univariate MARS analysis on (a)  $\text{CCS}_H$  and (b)  $\text{CCS}_{Na}$  models. blue lines are placed at 2% error bands and the purple at 5%. It is clear that the model still predicts well at higher values where there are less data but the prediction intervals are much larger to accommodate the uncertainty due to lack of data.

identification can be garnered by matching both of the CCS values observed with the predicted data.

The  $\text{CCS}_{Na}$  model herein presented also improves the prediction accuracy of the previously developed model by the authors.<sup>29</sup> In that work, we evaluated the performance of the ANN predictive model for sodium adducts, finding that deviations between predicted and experimental data were below 8.7% for 95% of the cases. However, the development of an exclusive model for the sodium adducts by MARS improves the prediction accuracy.

**3.1.3. Blind Testing of the Models.** Several reference standards were purchased from different research projects during the development of the predictors based on MARS. However, these newly available compounds were not included in the training and validation data sets used but were used to verify the utility of our prediction models for chemicals not previously considered in the training steps. Thus, model applicability can be extrapolated for upcoming RT and CCS predictions of real suspect compounds. Therefore, we calculated deviations between predicted and experimental data for this data set and compared the observed deviations with previously calculated accuracies at different percentiles (shown in Table 2). Table 3 depicts the experimental and predicted values of RT and CCS for the different adducts observed for the additional set of 25 reference standards. Moreover, the deviation between experimental and predicted is shown, and as it can be observed, the RT predictions are generally in agreement with the experimental data with the 95th percentile of the observed deviations ( $\pm 4.15$  min) being in the same range than that

observed during validation. However, *diphenyl hydrogen phosphate* showed a high deviation between experimental and predicted values. This can be potentially explained because of the lack of sufficient compounds featuring a P atom in their chemical structures in the initial training database. Hence, it is not surprising that for these molecular skeletons, the RT prediction does not fit precisely with the experimental data.

Furthermore, the vast majority of CCS values for  $[M + H]^+$  are in agreement with the values calculated using the  $\text{CCS}_H$  model. For these compounds, 95% of the cases showed deviations below  $\pm 3.71\%$ , yielding even better results than the initial database during model validation. Only *3,4-dichloroaniline* shows a deviation greater than 4%, which could be explained by the small CCS value calculated. When evaluating  $\text{CCS}_{Na}$ , higher deviations are observed concretely for the case of *di(2-ethylhexyl) terephthalate* and *vildagliptin* ( $-8.61$  and  $8.74\%$ , respectively). These deviations could be explained because of the particular chemical structures of the molecule such as the presence of an adamantyl group in *vildagliptin*, which has a large and rigid structure, or the high rotatability of alkyl chains in the *di(2-ethylhexyl) terephthalate*. However, if these adducts would be treated as outliers, 95% of the  $\text{CCS}_{Na}$  values show deviations of  $\pm 3.15\%$ , which is in accordance with the data obtained during method validation. Finally, for  $[M - H]^-$ , a small set of molecules was gathered, and all of them fit well within the  $\pm 5.8\%$  deviation.

**3.2. Open-Access Prediction Platform.** To aid future researchers working with UHPLC-IMS-HRMS, a free online webpage incorporating these models has been released. The models are available for the scientific community through [https://datascience-adelaideuniversity.shinyapps.io/Predicting\\_RT\\_and\\_CCS/](https://datascience-adelaideuniversity.shinyapps.io/Predicting_RT_and_CCS/). Figure 4 illustrates the layout of the online platform for the prediction of RT and CCS for both (de)protonated molecules or sodium adducts.

The operational of the platform is user-friendly and easy-to-follow. As an example, the step-by-step method to obtain prediction for omeprazole is shown. First, selection of which parameter is going to be predicted need to be done (Figure 4A). In this case, CCS for protonated molecules is selected by indicating “Select Response: Collision Cross Section” and “Sodiated: No”. After downloading the appropriate descriptors for the molecule of interest using Dragon v5.4 integrated within OChem ([www.ochem.eu](http://www.ochem.eu)),<sup>46</sup> those can be added in the corresponding editable fields (Figure 4B). The CCS value can, then, be predicted, and the output is shown together with their corresponding prediction intervals (Figure 4C). In this case, the CCS predicted value for the protonated molecule of omeprazole is  $181.51 \text{ \AA}^2$  with a prediction interval of  $171.93 - 190.08 \text{ \AA}^2$ . The experimental value for  $[M + H]^+$  for omeprazole is  $180.58 \text{ \AA}^2$ , denoting that the prediction only deviated by 0.52% from the experimental value.

The ease of prediction as well as the open access for this online platform is of great help for those researchers working on UHPLC-IMS-HRMS instruments who do not have an in-house-developed prediction model.

## 4. CONCLUSIONS

Three different prediction models using MARS have been developed for the prediction of RT, CCS for (de)protonated molecules, and CCS for sodium adducts. This is the first application of MARS for the prediction of RT and CCS data. In addition, the reported models are the first parallel prediction of RT and CCS data for the same instrument, facilitating the

**Table 3. Experimental and Predicted Values of RT and CCS for Additional Compounds Not Initially Included in Data Sets: Investigation of the Deviation of Predicted Values**

compound	retention time (min)			CCS <sub>H</sub> for [M + H] <sup>+</sup>			CCS <sub>Na</sub> for [M + Na] <sup>+</sup>			CCS <sub>H</sub> for [M - H] <sup>-</sup>		
	exp.	pred.	dev (min)	exp.	pred.	dev (%)	exp.	pred.	dev (%)	exp.	pred.	dev (%)
(-)-cotine	0.87	3.05	2.18	141.48	136.12	-3.79%						
3,4-dichloroaniline	7.92	6.21	-1.71	137.10	125.87	-8.19%						
3-hydroxyphenyl diphenyl phosphate	11.09	10.36	-0.73	174.31	178.94	2.66%	184.84	189.67	2.61%	180.89	178.94	-1.07%
5,6-dimethylbenzotriazole	6.74	4.38	-2.36	129.73	127.21	-1.94%				129.38	127.21	-1.68%
8-hydroxyquinoline	1.51	4.57	3.06	125.01	123.54	-1.18%						
amisulpride	2.46	1.99	-0.47	193.15	189.60	-1.84%						
antiblaze V6	10.84	14.49	3.65	208.45	207.32	-0.54%	208.45	212.12	1.76%			
benzotriazole	3.50	2.75	-0.75	121.49	117.94	-2.92%						
BCIPHP phosphate <sup>a</sup>	8.07	7.69	-0.38	159.22	157.05	-1.36%	165.14	166.24	0.67%			
caffeine	3.08	1.94	-1.14	136.62	136.37	-0.18%						
chlorotoluron	2.54	6.79	4.25	146.29	146.00	-0.20%	155.35	157.96	1.68%			
citalopram	6.49	5.21	-1.28	179.10	184.48	3.01%						
Di(2-ethylhexyl) terephthalate	16.86	15.02	-1.84				216.36	197.07.23	-8.61%			
diphenyl hydrogen phosphate	12.46	5.06	-7.41	152.45	151.61	-0.55%	161.58	162.65	0.66%	152.18	151.61	-0.38%
diphenylcresyl phosphate	7.36	11.07	3.72	175.28	178.08	1.60%						
metolachlor ESA <sup>b</sup>	7.89	4.99	-2.90	168.38	171.29	1.73%	175.57	179.13	2.03%	174.30	171.29	-1.73%
metoxuron	5.98	7.04	1.06	149.83	150.62	0.53%	158.51	161.17	1.68%			
mono(2-ethylhexyl) phthalate	12.73	11.75	-0.98							170.91	167.767215	-1.84%
monuron	6.68	5.67	-1.01	140.59	142.94	1.67%						
nicotine	0.69	1.11	0.42	138.34	134.77	-2.58%						
niflumic acid	11.51	10.86	-0.65	157.46	157.79	0.21%				156.92	157.79	0.55%
pirbuterol	1.30	1.28	-0.02	153.78	156.91	2.04%	160.02	165.52	3.44%			
prometon	6.74	7.50	0.76	156.67	155.56	-0.71%						
trietazine	10.81	8.91	-1.90	150.63	151.12	0.33%						
vildagliptin	1.38	1.67	0.29	176.98	174.62	-1.33%	172.29	187.35	8.74%			

<sup>a</sup>Bis(1-chloro-2-propyl) 1-hydroxy-2-propyl phosphate. <sup>b</sup>Metolachlor ethane sulfonic acid.

Predicting Retention Time or Collision Cross Section

Sodiated  
No

Select Response  
Collision Cross Section (CCS)

AMR  
93.813

Whetp  
1325.628

L1m  
19.402

PCR  
1.475

nRCHO  
0

LPRS  
114.103

MDOD  
18.938

Your selected response variable is Collision Cross Section (CCS)  
The best guess is a collision cross section of 181.51 with a 95% prediction interval of 172.93 to 190.08  
Compute

Example:  
Omeprazole

Empirical CCS for [M+H]<sup>+</sup>: 180.58 Å<sup>2</sup>  
Predicted CCS using CCS<sub>H</sub>: 181.51 Å<sup>2</sup> (0.52% ✓)

Web:  
[https://datascience-adelaideuniversity.shinyapps.io/Predicting\\_RT\\_and\\_CCS/](https://datascience-adelaideuniversity.shinyapps.io/Predicting_RT_and_CCS/)

**Figure 4.** Online platform for the prediction of RT and CCS data using univariate models. (A) Selection of response to predict, that is, RT, CCS for (de)protonated molecules or CCS for sodium adducts; (B) introduction of molecular descriptors for the molecule of interest; (C) output of the predictor model together with the prediction intervals. Example illustrated by omeprazole.

identification process of chemicals of emerging concern in SS and NTS strategies. The developed predictive models make use of a set of 26 molecular descriptors to predict RT and/or CCS values. The prediction accuracy achieved with these models bettered previously reported models in the literature by reducing

the deviation between predicted and experimental to  $\pm 2.32$  min for RT,  $\pm 4.05\%$  for CCS of protonated molecules,  $\pm 5.86\%$  for CCS of deprotonated molecules, and  $\pm 5.25\%$  for CCS of sodium adducts (95% confidence intervals). Additionally, a free access online platform has been released to enable the

application of these models to third-party laboratories interested in predicting RT and CCS data.

## DATA AND SOFTWARE AVAILABILITY

Data used for model development and validation are available in Table S1 of the Supporting Information as well as on the open-access online repository Zenodo (<https://zenodo.org/record/3966751#.Ymf5f9rP1aQ>, DOI: 10.5281/zenodo.3966751). Additionally, molecular descriptors are also shown in Table S1. Mathematical equations resulting from MARS model development are available for their implementation throughout the manuscript (eqs 1–3).

## ASSOCIATED CONTENT

### Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jcim.2c00847>.

Whole database of substances used for model development and validation including RT, CCS values, as well as the set of 1666 molecular descriptors (XLSX)

## AUTHOR INFORMATION

### Corresponding Authors

**Melissa Humphries** – School of Mathematical Sciences, University of Adelaide, SA-5005 Adelaide, Australia; Email: [melissa.humphries@adelaide.edu.au](mailto:melissa.humphries@adelaide.edu.au)

**Lubertus Bijlsma** – Environmental and Public Health Analytical Chemistry, Research Institute for Pesticides and Water, University Jaume I, E-12071 Castelló, Spain; [orcid.org/0000-0001-7005-8775](https://orcid.org/0000-0001-7005-8775); Email: [bijlsma@uji.es](mailto:bijlsma@uji.es)

### Authors

**Alberto Celma** – Environmental and Public Health Analytical Chemistry, Research Institute for Pesticides and Water, University Jaume I, E-12071 Castelló, Spain; Department of Aquatic Sciences and Assessment, Swedish University of Agricultural Sciences (SLU), SE-750 07 Uppsala, Sweden; [orcid.org/0000-0001-9763-8737](https://orcid.org/0000-0001-9763-8737)

**Richard Bade** – University of South Australia, Adelaide, UniSA: Clinical and Health Sciences, Health and Biomedical Innovation, Adelaide SA-5000 South Australia, Australia; Queensland Alliance for Environmental Health Sciences (QAEHS), The University of Queensland, Woolloongabba AUS-4102 Queensland, Australia; [orcid.org/0000-0003-2724-9183](https://orcid.org/0000-0003-2724-9183)

**Juan Vicente Sancho** – Environmental and Public Health Analytical Chemistry, Research Institute for Pesticides and Water, University Jaume I, E-12071 Castelló, Spain

**Félix Hernandez** – Environmental and Public Health Analytical Chemistry, Research Institute for Pesticides and Water, University Jaume I, E-12071 Castelló, Spain; [orcid.org/0000-0003-1268-3083](https://orcid.org/0000-0003-1268-3083)

Complete contact information is available at: <https://pubs.acs.org/10.1021/acs.jcim.2c00847>

### Author Contributions

<sup>#</sup>A.C. and R.B. are co-first authors.

### Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

A.C. acknowledges the Spanish Ministry of Economy and Competitiveness for his predoctoral grant (BES-2016-076914). L.B. acknowledges his fellowship funded by “la Caixa” Foundation. The project that gave rise to these results received the support of a fellowship from “la Caixa” Foundation (ID 10 010434). The fellowship code is LCF/BQ/PR21/11840012. Authors from University Jaume I acknowledge the financial support of the Spanish Ministry of Science, Innovation and Universities (RTI2018-097417-B-100), of Generalitat Valenciana (Research Group of Excellence Prometeo 2019/040), and of the University Jaume I of Castellón, Spain (project UJI-B2018-55 and UJI-B2020-19).

## REFERENCES

- (1) Hernández, F.; Bakker, J.; Bijlsma, L.; de Boer, J.; Botero-Coy, A. M.; Bruinen de Bruin, Y.; Fischer, S.; Hollender, J.; Kasprzyk-Hordern, B.; Lamoree, M.; López, F. J.; ter Laak, T. L.; van Leerdam, J. A.; Sancho, J. V.; Schymanski, E. L.; de Voogt, P.; Hogendoorn, E. A. The Role of Analytical Chemistry in Exposure Science: Focus on the Aquatic Environment. *Chemosphere* **2019**, *222*, 564–583.
- (2) Hollender, J.; van Bavel, B.; Dulio, V.; Farnen, E.; Furtmann, K.; Koschorreck, J.; Kunkel, U.; Krauss, M.; Munthe, J.; Schlabach, M.; Slobodnik, J.; Stroomberg, G.; Ternes, T.; Thomaidis, N. S.; Togola, A.; Tornero, V. High Resolution Mass Spectrometry-Based Non-Target Screening Can Support Regulatory Environmental Monitoring and Chemicals Management. *Environ. Sci. Eur.* **2019**, *31*, 42.
- (3) Schymanski, E. L.; Singer, H. P.; Slobodnik, J.; Ipolyi, I. M.; Oswald, P.; Krauss, M.; Schulze, T.; Haglund, P.; Letzel, T.; Grosse, S.; Thomaidis, N. S.; Bletsou, A.; Zwiener, C.; Ibáñez, M.; Portolés, T.; De Boer, R.; Reid, M. J.; Onghena, M.; Kunkel, U.; Schulz, W.; Guillon, A.; Noyon, N.; Leroy, G.; Bados, P.; Bogianni, S.; Stipanichev, D.; Rostkowski, P.; Hollender, J. Non-Target Screening with High-Resolution Mass Spectrometry: Critical Review Using a Collaborative Trial on Water Analysis. *Anal. Bioanal. Chem.* **2015**, *407*, 6237–6255.
- (4) Samanipour, S.; Langford, K.; Reid, M. J.; Thomas, K. V. A Two Stage Algorithm for Target and Suspect Analysis of Produced Water via Gas Chromatography Coupled with High Resolution Time of Flight Mass Spectrometry. *J. Chromatogr. A* **2016**, *1463*, 153–161.
- (5) Alygizakis, N. A.; Oswald, P.; Thomaidis, N. S.; Schymanski, E. L.; Aalizadeh, R.; Schulze, T.; Oswaldova, M.; Slobodnik, J. NORMAN Digital Sample Freezing Platform: A European Virtual Platform to Exchange Liquid Chromatography High Resolution-Mass Spectrometry Data and Screen Suspects in “Digitally Frozen” Environmental Samples. *TrAC, Trends Anal. Chem.* **2019**, *115*, 129–137.
- (6) Aalizadeh, R.; Nika, M.-C.; Thomaidis, N. S. Development and Application of Retention Time Prediction Models in the Suspect and Non-Target Screening of Emerging Contaminants. *J. Hazard. Mater.* **2019**, *363*, 277–285.
- (7) Bijlsma, L.; Berntssen, M. H. G.; Merel, S. A Refined Nontarget Workflow for the Investigation of Metabolites through the Prioritization by in Silico Prediction Tools. *Anal. Chem.* **2019**, *91*, 6321–6328.
- (8) Schymanski, E. L.; Jeon, J.; Gulde, R.; Fenner, K.; Ruff, M.; Singer, H. P.; Hollender, J. Identifying Small Molecules via High Resolution Mass Spectrometry: Communicating Confidence. *Environ. Sci. Technol.* **2014**, *48*, 2097–2098.
- (9) Celma, A.; Sancho, J. V.; Schymanski, E. L.; Fabregat-Safont, D.; Ibáñez, M.; Goshawk, J.; Barknowitz, G.; Hernández, F.; Bijlsma, L. Improving Target and Suspect Screening High-Resolution Mass Spectrometry Workflows in Environmental Analysis by Ion Mobility Separation. *Environ. Sci. Technol.* **2020**, *54*, 15120–15131.
- (10) Regueiro, J.; Negreira, N.; Berntssen, M. H. G. Ion-Mobility-Derived Collision Cross Section as an Additional Identification Point for Multiresidue Screening of Pesticides in Fish Feed. *Anal. Chem.* **2016**, *88*, 11169–11177.
- (11) McCullagh, M.; Giles, K.; Richardson, K.; Stead, S.; Palmer, M. Investigations into the Performance of Travelling Wave Enabled



Conventional and Cyclic Ion Mobility Systems to Characterise Protomers of Fluoroquinolone Antibiotic Residues. *Rapid Commun. Mass Spectrom.* **2019**, *33*, 11–21.

(12) Celma, A.; Ahrens, L.; Gago-Ferrero, P.; Hernández, F.; López, F.; Lundqvist, J.; Pitarch, E.; Sancho, J. V.; Wiberg, K.; Bijlsma, L. The Relevant Role of Ion Mobility Separation in LC-HRMS Based Screening Strategies for Contaminants of Emerging Concern in the Aquatic Environment. *Chemosphere* **2021**, *280*, No. 130799.

(13) Mollerup, C. B.; Mardal, M.; Dalsgaard, P. W.; Linnet, K.; Barron, L. P. Prediction of Collision Cross Section and Retention Time for Broad Scope Screening in Gradient Reversed-Phase Liquid Chromatography-Ion Mobility-High Resolution Accurate Mass Spectrometry. *J. Chromatogr. A* **2018**, *1542*, 82–88.

(14) Gabelica, V.; Marklund, E. Fundamentals of Ion Mobility Spectrometry. *Curr. Opin. Chem. Biol.* **2018**, *42*, 51–59.

(15) Lee, J. W. Basics of Ion Mobility Mass Spectrometry. *Mass Spectrom. Lett.* **2017**, *8*, 79–89.

(16) Aalizadeh, R.; Alygizakis, N. A.; Schymanski, E. L.; Krauss, M.; Schulze, T.; Ibáñez, M.; McEachran, A. D.; Chao, A.; Williams, A. J.; Gago-Ferrero, P.; Covaci, A.; Moschet, C.; Young, T. M.; Hollender, J.; Slobodnik, J.; Thomaidis, N. S. Development and Application of Liquid Chromatographic Retention Time Indices in HRMS-Based Suspect and Nontarget Screening. *Anal. Chem.* **2021**, *93*, 11601–11611.

(17) Stoffel, R.; Quilliam, M. A.; Hardt, N.; Fridstrom, A.; Witting, M. N-Alkylpyridinium Sulfonates for Retention Time Indexing in Reversed-Phase-Liquid Chromatography-Mass Spectrometry-Based Metabolomics. *Anal. Bioanal. Chem.* **2022**, *414*, 7387.

(18) Celma, A.; Bijlsma, L.; López, F. J.; Sancho, J. V. Development of a Retention Time Interpolation Scale (RTi) for Liquid Chromatography Coupled to Mass Spectrometry in Both Positive and Negative Ionization Modes. *J. Chromatogr. A* **2018**, *1568*, 101–107.

(19) Yeung, D.; Klaassen, N.; Mizero, B.; Spicer, V.; Krokhin, O. V. Peptide Retention Time Prediction in Hydrophilic Interaction Liquid Chromatography: Zwitter-Ionic Sulfoalkylbetaine and Phosphorylcholine Stationary Phases. *J. Chromatogr. A* **2020**, *1619*, No. 460909.

(20) Kensert, A.; Bouwmeester, R.; Efthymiadis, K.; Van Broeck, P.; Desmet, G.; Cabooter, D. Graph Convolutional Networks for Improved Prediction and Interpretability of Chromatographic Retention Data. *Anal. Chem.* **2021**, *93*, 15633–15641.

(21) Yang, J. J.; Han, Y.; Mah, C. H.; Wanjaya, E.; Peng, B.; Xu, T. F.; Liu, M.; Huan, T.; Fang, M. L. Streamlined MRM Method Transfer between Instruments Assisted with HRMS Matching and Retention-Time Prediction. *Anal. Chim. Acta* **2020**, *1100*, 88–96.

(22) Barron, L. P.; McEneff, G. L. Gradient Liquid Chromatographic Retention Time Prediction for Suspect Screening Applications: A Critical Assessment of a Generalised Artificial Neural Network-Based Approach across 10 Multi-Residue Reversed-Phase Analytical Methods. *Talanta* **2016**, *147*, 261–270.

(23) Stanstrup, J.; Neumann, S.; Vrhovšek, U. PredRet: Prediction of Retention Time by Direct Mapping between Multiple Chromatographic Systems. *Anal. Chem.* **2015**, *87*, 9421–9428.

(24) Bade, R.; Bijlsma, L.; Sancho, J. V.; Hernández, F. Critical Evaluation of a Simple Retention Time Predictor Based on LogKow as a Complementary Tool in the Identification of Emerging Contaminants in Water. *Talanta* **2015**, *139*, 143–149.

(25) Bade, R.; Bijlsma, L.; Miller, T. H.; Barron, L. P.; Sancho, J. V.; Hernández, F. Suspect Screening of Large Numbers of Emerging Contaminants in Environmental Waters Using Artificial Neural Networks for Chromatographic Retention Time Prediction and High Resolution Mass Spectrometry Data Analysis. *Sci. Total Environ.* **2015**, *538*, 934–941.

(26) Bonini, P.; Kind, T.; Tsugawa, H.; Barupal, D. K.; Fiehn, O. Retip: Retention Time Prediction for Compound Annotation in Untargeted Metabolomics. *Anal. Chem.* **2020**, *92*, 7515–7522.

(27) Domingo-Almenara, X.; Guijas, C.; Billings, E.; Montenegro-Burke, J. R.; Uritboonthai, W.; Aisporna, A. E.; Chen, E.; Benton, H. P.; Siuzdak, G. The METLIN Small Molecule Dataset for Machine Learning-Based Retention Time Prediction. *Nat. Commun.* **2019**, *10*, 5811.

(28) Yang, Q.; Ji, H.; Fan, X.; Zhang, Z.; Lu, H. Retention Time Prediction in Hydrophilic Interaction Liquid Chromatography with Graph Neural Network and Transfer Learning. *J. Chromatogr. A* **2021**, *1656*, No. 462536.

(29) Bijlsma, L.; Bade, R.; Celma, A.; Mullin, L.; Cleland, G.; Stead, S.; Hernandez, F.; Sancho, J. V. Prediction of Collision Cross-Section Values for Small Molecules: Application to Pesticide Residue Analysis. *Anal. Chem.* **2017**, *89*, 6583–6589.

(30) Zhou, Z.; Tu, J.; Xiong, X.; Shen, X.; Zhu, Z.-J. LipidCCS: Prediction of Collision Cross-Section Values for Lipids with High Precision To Support Ion Mobility–Mass Spectrometry-Based Lipidomics. *Anal. Chem.* **2017**, *89*, 9559–9566.

(31) Plante, P.-L.; Francovic-Fontaine, É.; May, J. C.; McLean, J. A.; Baker, E. S.; Lavolette, F.; Marchand, M.; Corbeil, J. Predicting Ion Mobility Collision Cross-Sections Using a Deep Neural Network: DeepCCS. *Anal. Chem.* **2019**, *91*, 5191.

(32) Ross, D. H.; Cho, J. H.; Xu, L. Breaking Down Structural Diversity for Comprehensive Prediction of Ion-Neutral Collision Cross Sections. *Anal. Chem.* **2020**, *92*, 4548.

(33) Zhou, Z.; Shen, X.; Tu, J.; Zhu, Z.-J. Large-Scale Prediction of Collision Cross-Section Values for Metabolites in Ion Mobility-Mass Spectrometry. *Anal. Chem.* **2016**, *88*, 11084–11091.

(34) Colby, S. M.; Thomas, D. G.; Nunez, J. R.; Baxter, D. J.; Glaesemann, K. R.; Brown, J. M.; Pirrung, M. A.; Govind, N.; Teegarden, J. G.; Metz, T. O.; Renslow, R. S. ISiCLE: A Quantum Chemistry Pipeline for Establishing in Silico Collision Cross Section Libraries. *Anal. Chem.* **2019**, *91*, 4346–4356.

(35) Gonzales, G. B.; Smagghe, G.; Coelus, S.; Adriaenssens, D.; De Winter, K.; Desmet, T.; Raes, K.; Van Camp, J. Collision Cross Section Prediction of Deprotonated Phenolics in a Travelling-Wave Ion Mobility Spectrometer Using Molecular Descriptors and Chemometrics. *Anal. Chim. Acta* **2016**, *924*, 68–76.

(36) Zhou, Z.; Luo, M.; Chen, X.; Yin, Y.; Xiong, X.; Wang, R.; Zhu, Z. Ion Mobility Collision Cross-Section Atlas for Known and Unknown Metabolite Annotation in Untargeted Metabolomics. *Nat. Commun.* **2020**, *11*, 4334.

(37) Ewing, S. A.; Donor, M. T.; Wilson, J. W.; Prell, J. S. Collidoscope: An Improved Tool for Computing Collisional Cross-Sections with the Trajectory Method. *J. Am. Soc. Mass Spectrom.* **2017**, *28*, 587–596.

(38) Lee, J. W.; Lee, H. H. L.; Davidson, K. L.; Bush, M. F.; Kim, H. I. Structural Characterization of Small Molecular Ions by Ion Mobility Mass Spectrometry in Nitrogen Drift Gas: Improving the Accuracy of Trajectory Method Calculations. *Analyst* **2018**, *143*, 1786–1796.

(39) Zanutto, L.; Heerdt, G.; Souza, P. C. T.; Araujo, G.; Skaf, M. S. High Performance Collision Cross Section Calculation-HPCCS. *J. Comput. Chem.* **2018**, *39*, 1675–1681.

(40) Falchi, F.; Bertozzi, S. M.; Ottonello, G.; Ruda, G. F.; Colombano, G.; Fiorelli, C.; Martucci, C.; Bertorelli, R.; Scarpelli, R.; Cavalli, A.; Bandiera, T.; Armirotti, A. Kernel-Based, Partial Least Squares Quantitative Structure-Retention Relationship Model for UPLC Retention Time Prediction: A Useful Tool for Metabolite Identification. *Anal. Chem.* **2016**, *88*, 9510–9517.

(41) Put, R.; Xu, Q. S.; Massart, D. L.; Vander Heyden, Y. Multivariate Adaptive Regression Splines (MARS) in Chromatographic Quantitative Structure–Retention Relationship Studies. *J. Chromatogr. A* **2004**, *1055*, 11–19.

(42) Souihi, A.; Mohai, M. P.; Palm, E.; Malm, L.; Krueve, A. MultiConditionRT: Predicting Liquid Chromatography Retention Time for Emerging Contaminants for a Wide Range of Eluent Compositions and Stationary Phases. *J. Chromatogr. A* **2022**, *1666*, No. 462867.

(43) Friedman, J. H. Multivariate Adaptive Regression Splines. *Ann. Stat.* **1991**, *19*, 1–67.

(44) Put, R.; Vander Heyden, Y. Review on Modelling Aspects in Reversed-Phase Liquid Chromatographic Quantitative Structure–Retention Relationships. *Anal. Chim. Acta* **2007**, *602*, 164–172.

(45) Celma, A.; Fabregat-Safont, D.; Ibáñez, M.; Bijlsma, L.; Hernández, F.; Sancho, J. V. S61 | UJICCSLIB | Collision Cross

Section (CCS) Library from UJI (Version NORMAN-SLE S61.0.1.2) [Data Set]. *Zenodo* 2019, DOI: 10.5281/zenodo.3966751.

(46) Sushko, I.; Novotarskyi, S.; Körner, R.; Pandey, A. K.; Rupp, M.; Teetz, W.; Brandmaier, S.; Abdelaziz, A.; Prokopenko, V. V.; Tanchuk, V. Y.; Todeschini, R.; Varnek, A.; Marcou, G.; Ertl, P.; Potemkin, V.; Grishina, M.; Gasteiger, J.; Schwab, C.; Baskin, I. I.; Palyulin, V. A.; Radchenko, E. V.; Welsh, W. J.; Kholodovych, V.; Chekmarev, D.; Cherkasov, A.; Aires-de-Sousa, J.; Zhang, Q.-Y.; Bender, A.; Nigsch, F.; Patiny, L.; Williams, A.; Tkachenko, V.; Tetko, I. V. Online Chemical Modeling Environment (OCHEM): Web Platform for Data Storage, Model Development and Publishing of Chemical Information. *J. Comput.-Aided Mol. Des.* 2011, 25, 533–554.

(47) Team, R. C. R: *A Language and Environment for Statistical Computing*; R Foundatin for Statistical Computing: Vienna, Austria2020.

(48) Milborrow, S.. *Derived from mda.mars by T. Hastie and R. Tibshirani. Earth: Multivariate Adaptive Regression Splines (R Package v5.3.0)*, 2011.

(49) Viswanadhan, V. N.; Ghose, A. K.; Revankar, G. R.; Robins, R. K. Atomic Physicochemical Parameters for Three Dimensional Structure Directed Quantitative Structure-Activity Relationships. 4. Additional Parameters for Hydrophobic and Dispersive Interactions and Their Application for an Automated Superposition of Certain. *J. Chem. Inf. Model.* 1989, 29, 163–172.

(50) Tetko, I. V.; Tanchuk, V. Y. Application of Associative Neural Networks for Prediction of Lipophilicity in ALOGPS 2.1 Program. *J. Chem. Inf. Comput. Sci.* 2002, 42, 1136–1145.

(51) Tetko, I. V.; Gasteiger, J.; Todeschini, R.; Mauri, A.; Livingstone, D.; Ertl, P.; Palyulin, V. A.; Radchenko, E. V.; Zefirov, N. S.; Makarenko, A. S.; Tanchuk, V. Y.; Prokopenko, V. V. Virtual Computational Chemistry Laboratory–Design and Description. *J. Comput.-Aided Mol. Des.* 2005, 19, 453–463.