

# Modeling the ionization efficiency of small molecules in positive electrospray ionization

Dimitri Abrahamsson<sup>1\*</sup>, Lelouda-Athanasia Koronaiou<sup>2,3</sup>, Trevor Johnson<sup>1</sup>, Dimitra A. Lambropoulou<sup>2,3</sup>

<sup>1</sup> Department of Pediatrics, New York University Grossman School of Medicine, New York, New York 10016, United States

<sup>2</sup> Laboratory of Environmental Pollution Control, Department of Chemistry, Aristotle University of Thessaloniki, University Campus 54124 Thessaloniki, Greece

<sup>3</sup> Center for Interdisciplinary Research and Innovation (CIRI-AUTH), Balkan Center, Thessaloniki 57001, Greece

\*Corresponding author: Dimitri Abrahamsson, [dimitri.abrahamsson@gmail.com](mailto:dimitri.abrahamsson@gmail.com)

Keywords: electrospray ionization, predictive modeling, relative response factor, molecular dynamics simulations

## ABSTRACT

Technological advancements in liquid chromatography (LC) electrospray ionization (ESI) high resolution mass spectrometry (HRMS) have made it an increasingly popular analytical technique in non-targeted analysis (NTA) of environmental and biological samples. One critical limitation of current methods in NTA is the lack of available analytical standards for many of the compounds detected in biological and environmental samples. Computational approaches can provide estimates of concentrations by modeling the ionization efficiency of a compound expressed as the relative response factor (RRF). In this paper, we explore the application of molecular dynamics (MD) in the development of a computational workflow for predicting RRF. We obtained measurements of RRF for 48 compounds with LC - quadrupole time-of-flight (QTOF) MS and calculated their RRF by dividing the observed peak areas by their concentrations. We used the CGenFF force field to generate the topologies and GROMACS to conduct the (MD) simulations ( $t = 1$  ns). We calculated the Lennard-Jones and Coulomb interactions between the analytes and all other molecules in the ESI droplet, which were then used to construct a multilinear regression model for predicting RRF. The best performing model showed a coefficient of determination ( $R^2$ ) of 0.82 and a mean absolute error (MAE) of 0.13 log units. This performance is comparable to other predictive models including machine learning models. While there is a need for further evaluation of diverse chemical structures, our approach showed promise in predictions of RRF.

## 46 INTRODUCTION

47 Recent technological breakthroughs in high-resolution mass spectrometry (HRMS) have made it an  
48 increasingly popular technology in environmental analysis<sup>1</sup>, biomonitoring<sup>2,3</sup> and metabolomics<sup>4</sup>. Often  
49 combined with suspect screening (SS) and non-targeted analysis (NTA) workflows, HRMS has shown  
50 great promise in the discovery of lesser-known chemical structures and in comprehensively  
51 characterizing the chemical composition of complex mixtures.<sup>2,3,5</sup> One critical challenge associated with  
52 the application of HRMS in SS or NTA is the limited availability of analytical standards for many  
53 anthropogenic / synthetic chemicals and many endogenously produced metabolites.<sup>6,7</sup>

54 When it comes to endogenous metabolites, the lack of commercially available analytical standards is  
55 often due to certain compounds being less well-characterized and thus not yet synthesized or purified.  
56 For anthropogenic chemicals, one of the reasons is that chemical manufacturers in the U.S. are not  
57 required to produce analytical standards for the chemicals that they manufacture and release to the  
58 environment.<sup>7</sup> One exception to this rule is pesticides.<sup>7</sup> It is important to note at this stage that the  
59 requirement for analytical standards does not extend to the transformation and breakdown products of  
60 these chemicals. So even in a hypothetical scenario where manufacturers would be required to  
61 produce analytical standards, that would cover only the parent compounds and not all the  
62 transformation products. The U.S. Environmental Protection Agency (EPA) has prioritized about 1.2  
63 million chemicals of environmental importance and has created a database called EPA's CompTox  
64 Chemicals Dashboard (henceforth referred to as "the dashboard").<sup>8</sup> Nuñez et al.<sup>6</sup> estimated that out of  
65 the 1.2 million chemicals on EPA's Dashboard, less than 2% are available as analytical standards.

66 There is thus a need to develop computational approaches to confirm and quantify detected  
67 compounds without analytical standards.<sup>6,9,10</sup> While detection and tentative confirmation can be  
68 achieved with MS/MS libraries, quantification remains a more challenging task.<sup>11,12</sup> Liquid

Chromatography (LC) – Electrospray Ionization (ESI) HRMS is one of the most commonly used HRMS techniques in SS and NTA studies. One critical challenge in ESI is that abundances expressed as peak areas or peak heights are not easily translatable to concentrations. Two compounds of the same concentration can exhibit peak areas that differ by 3 orders of magnitude because of differences in ionization<sup>12,13</sup>. Abundances may be used as surrogate for concentrations in certain situations when comparing the same chemical across different samples, however, they cannot be used to compare two chemicals to each other.<sup>14</sup>

While ESI is extensively used in mass spectrometry for the analysis of both small (e.g., metabolites and large molecules (e.g., proteins), the precise mechanism has not been fully understood. Briefly, during ESI, the solution containing the analyte passes through a metal capillary that is charged at an electric potential of thousands of Volts (kV). The solution forms a tip at the end of the capillary known as a Taylor cone that emits a spray of fine droplets. The droplets start in the  $\mu\text{m}$  range and shrink in size as they undergo evaporation often accelerated with heating of the capillary. The density of the charged ions in the droplet is controlled by repulsive Coulombic forces between positively charged ions. The upper limit of that density is described as the Rayleigh stability limit:

$$z_R = \frac{8\pi}{e} \sqrt{\epsilon_0 \gamma R^3} \quad (\text{eq. 1})$$

where,  $e$  is the number of elementary charged particles,  $\epsilon_0$  is the vacuum permittivity,  $\gamma$  is the surface tension of the droplet and  $R$  is the droplet radius.

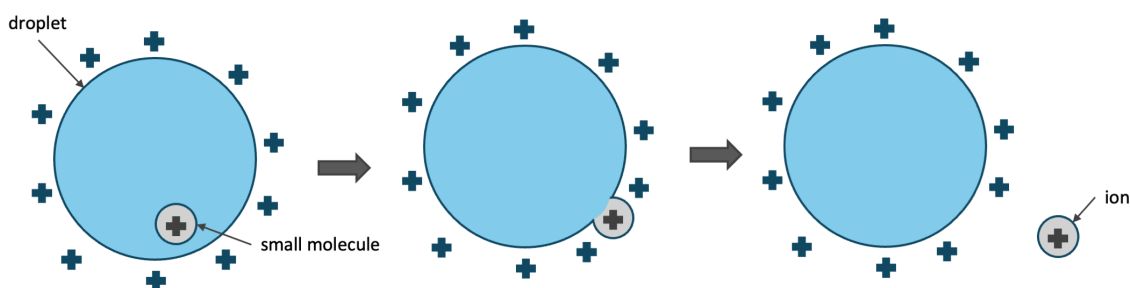
Conceptual models have been proposed to describe the process that molecules undergo to become ionized and transferred to the gas phase during ESI. We focus our discussion on the commonly used positive ESI under which positive ions are formed. Small molecules (< 1000 Da) are thought to ionize and be transferred to the gas phase by the ion evaporation model (IEM). According to IEM, the analyte is protonated already while inside the droplet and eventually moves from the center towards the

92 surface of the droplet. As the positively charged analyte meets the positively charged solvent molecules  
93 on the surface of the droplet, the ion is transferred to the gas phase through repulsive forces of  
94 positively charged ions and by the excess droplet charge.

95 Larger molecules such as globular proteins (natively folded proteins) are thought to ionize and be  
96 transferred to the gas phase by the Charged Residue Model (CRM). According to CRM, solvent droplets  
97 containing a single protein molecule gradually evaporate to dryness and as the solvent molecules  
98 evaporate, the charge is transferred to the analyte. The droplets remain close the Rayleigh stability limit  
99 while evaporating, which indicates that the droplet loses some of the electric charge as it shrinks in  
100 size. This is supported by experimental observations where the size of solvent droplets positively  
101 correlated with the droplet charge following an exponential curve. Contrary to small molecules, the  
102 ejection of globular proteins is not kinetically favorable. The repulsive forces of the excess surface  
103 charge are not sufficiently strong for the molecule to be ejected and transferred to the gas phase. CRM  
104 is supported by experimental evidence where ionization of globular proteins produces ions with a  
105 charge of  $[M + z_R H]^{z_R+}$ , where  $z_R$  is the charge at the Rayleigh limit of water droplets of the same size  
106 as the globular protein.

107 While these two models are considered distinct, both of these two processes could apply to small or  
108 medium size molecules, especially in the case of heated ESI where the evaporation of the water  
109 droplets is assisted through heating of the capillary. This is also supported by molecular dynamics (MD)  
110 simulations studies where native (unmodified) carbohydrates ionize through CRM, while their  
111 permethylated derivatives ionize through IEM.<sup>15</sup>

### Ion Evaporation Model (IEM)



### Charged Residue Model (CRM)

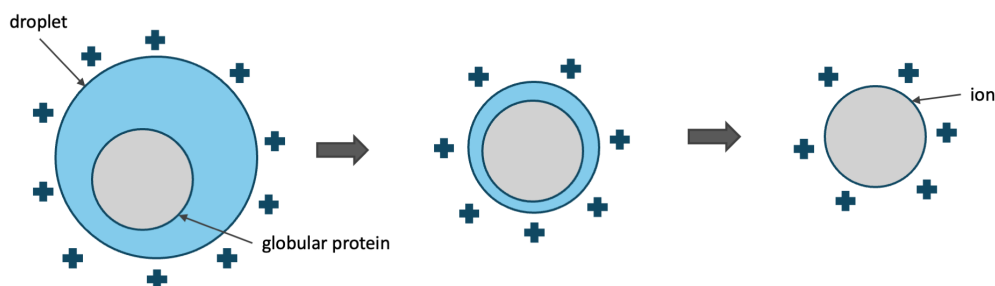


Figure 1: Conceptual models describing the mechanism of electrospray ionization. Small molecules are thought to be ionized through the ion evaporation model (IEM), while larger molecules such as globular proteins are thought to ionize through the charged residue model (CRM).

Ionization efficiency is a property that has proven to be difficult to predict from a small number of chemical descriptors or physicochemical properties. Numerous investigations have examined the correlation between ionization efficiencies and diverse physicochemical properties, including but not limited to their  $pK_a$ ,  $\log P$ , surface area, charge delocalization, and gas-phase proton affinity.<sup>16–21</sup> These findings have led to the development of various predictive models for ionization efficiencies<sup>19,22–25</sup>, utilizing both the physicochemical properties of analytes and solvent characteristics as fundamental inputs. These models commonly rely on parameters associated with the analyte's hydrophobicity (e.g.,  $\log P$ , WAPS, WANS, C/H ratio) and ionizability (such as  $pK_a$  and the degree of ionization)<sup>26</sup>. Data-driven approaches involving machine learning, such as random forest models, have shown great promise in

127 providing predictions with reasonable uncertainties.<sup>11,12,27–29</sup> It should be noted, however, that these  
128 approaches require large datasets in the order of 100s of chemicals with diverse structures and  
129 properties and the predictions are tied to a specific method and instrumentation.

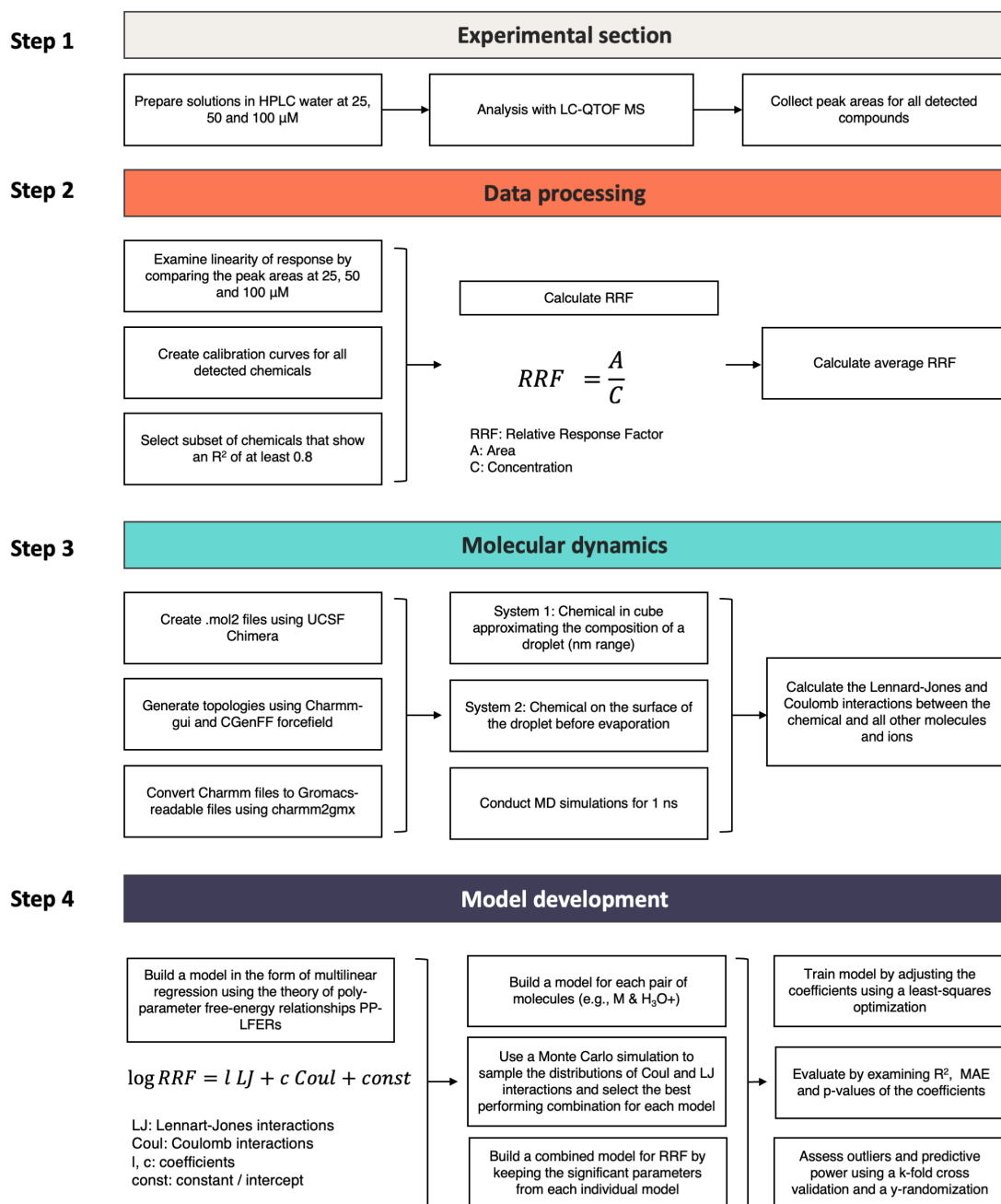
130 Theory-driven approaches that are based on quantum chemistry and computational chemistry  
131 principles could provide an alternative when large datasets are not available or are difficult to obtain.  
132 Molecular dynamic simulations (MD) have been previously applied in a limited number of studies to  
133 understand the mechanism of ionization in salt ions, peptides and proteins.<sup>30–34</sup> However, to the best of  
134 our knowledge, there appears to be little on the study of small molecules and especially in predictions  
135 of their ionization efficiency using MD simulations. Our study aims to fill this gap by employing  
136 molecular dynamics to model the behavior of chemicals in the ionization chamber and evaluate the  
137 potential of such theory-driven approaches to make predictions and assess their uncertainties.

## 150 MATERIALS AND METHODS

### 151 Workflow diagram

152 The individual steps of the experimental and computational aspects of the study are presented in

153 Figure 2.



154

155 Figure 2: Workflow diagram for the processing steps in the experimental and computational parts of



156 the study

157 Experimental section

158 Chemicals and solutions

159 The analytes were provided by the US EPA for the purposes of this study and were developed as part  
160 of EPA's Non-Targeted Analysis Collaborative Trial (ENTACT). The preparation of the chemical mixtures is  
161 described in detail in the study of Ulrich et al.<sup>35</sup> For the purposes of this study, we used mixtures 504,  
162 506 and 508. The mixtures were diluted in a series of dilutions first with methanol (99.9% Millipore  
163 Sigma) and then with HPLC water (99.9% Millipore Sigma) from 20 mM to 100, 50 and 25  $\mu$ M with a  
164 final water content of > 99%. For the purposes of this study, we selected chemicals that ionize in  
165 positive electrospray ionization mode (ESI+) and whose calibration curves showed an  $R^2$  of 0.8 or  
166 higher. The chemical structures and chemical identifiers of the chemicals involved in the study (n=48)  
167 are shown in Supplemental Spreadsheet 1. The complete chemical lists from mixtures 504, 506 and 508  
168 are also shown in Supplemental Spreadsheet 1. We should note at this point that characterizing the  
169 mixtures and maximizing the number of detected and identified compounds is beyond the scope of the  
170 study. As our study requires extensive computations that take days to complete, we have to limit our  
171 efforts to a small subset of compounds that satisfies the criteria of detection (< 5 ppm mass difference  
172 from the monoisotopic mass) and linearity ( $R^2 \geq 0.8$ ).

173

174 Instrumental analysis

175 The mixtures were analyzed with an Agilent 1290 ultra-high performance liquid chromatography  
176 (UPLC) coupled to an Agilent quadrupole time-of-flight (QTOF) mass spectrometer. The UPLC was  
177 equipped with an Agilent Eclipse Plus C18 column (2.1 x 100 mm, 1.8  $\mu$ M) for the chromatographic  
178 separation of the analytes. The mobile phase consisted of the two following solutions. Solution A: HPLC

179 water (Sigma Aldrich,  $\geq 99.5\%$ ) with 0.1% methanol (Sigma Aldrich, 99.9%) and 5 mM ammonium  
180 acetate (Sigma Aldrich,  $\geq 98\%$ ). Solution B: 90 %methanol with 10% HPLC water and 5 mM ammonium  
181 acetate. The two solutions were mixed under the following gradient program: 0 min 10% B and 90% A,  
182 0– 15 min increase to 100% B, 16–20 min equilibration at 100% B. The solvent gradient over time is also  
183 shown in Figure S1. All mixtures were injected twice at an injection volume of 5  $\mu\text{L}$ . Two no-injection  
184 blanks and one HPLC water blank were also analyzed in the beginning of the sequence.

185 The instrument was operated in both positive electrospray ionization mode (ESI+) and full scan mass  
186 spectra (MS1) were acquired in the range of 100-1000 Da with a resolving power of 40,000 and a mass  
187 accuracy of  $<5$  ppm. The instrument was calibrated before the analysis and the mass difference was  
188 corrected with reference standards using masses 121.050873 and 922.009798 for positive ionization  
189 mode.

190

#### 191 Data collection and file processing

192 All the collected data files were processed with MS-DIAL, an open-source software for mass  
193 spectrometry that was developed by the University of California, Davis and by RIKEN (Japan). The  
194 detected features were aligned across samples and were matched to the monoisotopic masses of the  
195 chemicals contained in the mixtures within a 10-ppm mass difference (Supplementary Spreadsheet).  
196 The peak areas of the analytes were calculated by taking the average of the duplicate injections and  
197 they were corrected by subtracting the average area measured in the blanks. The MS-DIAL settings and  
198 parameters used to process the data files are presented in Supplementary Spreadsheet.

199

#### 200 Calculation of RRF

201 Ionization efficiency describes the extent to which molecules of an analyte in the liquid phase can  
202 transition to the gas phase as ions during the process of electrospray ionization. The ionization  
203 efficiency of an analyte can be mathematically described by the relative response factor of the analyte  
204 as follows:

$$205 \quad RRF = \frac{A}{C} \quad eq. 2$$

206 Where, RRF is the relative response factor, A is the abundance (peak area or peak height) and C is the  
207 concentration of the analyte.

208

209 Molecular dynamics simulations

210 Input generation

211 We generated mol2 format files for all chemicals in the dataset using UCSF Chimera and using SMILES  
212 as inputs for the mol2 files. The protonation of each molecule was determined by generating pKa  
213 diagrams for each chemical using Chemaxon and Chemicalize<sup>36</sup> and identifying the dominant species at  
214 the pH that would be relevant from our experiments (pH=5).<sup>37</sup> All the pKa diagrams and the  
215 protonation states are uploaded as .png images on GitHub under  
216 <https://github.com/dimitriabrahamsson/electro-spray> .

217

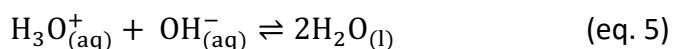
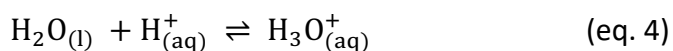
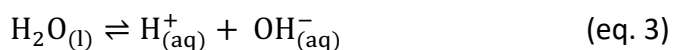
218 Topology generation

219 Topologies were generated with CGenFF force field (CHARMM General Force Field) using the  
220 CHARMM-GUI<sup>38</sup> online platform (<https://www.charmm-gui.org/>) and the mol2 files from the previous  
221 step. The protonation of the analytes was examined once more to ensure that it was correct and that  
222 no changes were made while importing the mol2 files in GHARMM-GUI. The generated files from

CHARMM-GUI were then converted to GROMACS format using the charmm2gromacs-pvm.py script (uploaded on GitHub under <https://github.com/dimitriabrahamsson/electro-spray> ).

## System preparation

All preparation steps were conducted using GROMACS (version 2023.2). GROMACS uses periodic boundary conditions (PBC) where the atoms of the simulation system are put into a space-filling box, which is surrounded by translated copies of itself. Thus, the system does not have finite borders during the simulation, but it allows for the removal of PBC for post-simulation calculations.<sup>39</sup> Two different systems were considered for the MD simulations. System 1 aimed at approximating the composition of a nm electrospray droplet which included the analyte, the H<sup>+</sup> and OH<sup>-</sup> ions produced during hydrolysis, and the water and methanol molecules as shown in the first step of IEM in Figure 1. With the second system (System 2), we aimed to approximate the composition of the droplet surface and the analyte located at the surface before evaporation as shown in the second step of IEM in Figure 1. One critical challenge when describing both systems is describing the H<sup>+</sup> ions in solution. In water, the H<sup>+</sup> ions, also referred to as H<sub>(aq)</sub><sup>+</sup>, produced from the hydrolysis of water molecules react with other water molecules to form hydronium, H<sub>3</sub>O<sub>(aq)</sub><sup>+</sup>, also known as oxonium, following the reactions below:



## 246 System 1

247 The molecule of hydronium was described using a TIP3P that included an additional hydrogen (3 H in  
248 total) and was modified to include the specific parameters for  $\text{H}_3\text{O}^+$  from the study of Wolf and  
249 Groenhof.<sup>40</sup> The distance between O and H ( $r_{\text{OH}}$ ) was set at  $1.02\text{\AA}$ , the angle for H-O-H ( $\theta_{\text{HOH}}$ ) was set at  
250  $112^\circ$ , and the charges for O ( $q_{\text{O}}$ ) and H ( $q_{\text{H}}$ ) were set at  $-0.59$  and  $0.53e$ . The droplet was represented by  
251 a three-dimensional cube and the size was set at  $64\text{ nm}^3$  ( $4 \times 4 \times 4\text{ nm}$ ). The number of  $\text{H}_3\text{O}^+$  molecules  
252 was approximated based the experimental observations of Smith et al.<sup>41</sup> who determined the charge  
253 (e) of water and methanol droplets in ESI+ as a function of the droplet diameter. The calculations are  
254 described in detail in Text S1 in SI.

255 The numbers of water (TIP3P) and methanol molecules were determined based on the gradient  
256 mixing (Figure S1) of the two solvents during LC and based on the retention time of each chemical  
257 (Supplemental Spreadsheet). This means that for every chemical the number of water and methanol  
258 molecules was different depending on when it was eluted from the LC column. Previous MD studies<sup>42</sup>  
259 on ESI droplets have also suggested that the amount of methanol in the droplet plays a critical role in  
260 the ionization efficiency of the analytes. As the volume of methanol increases, the evaporative rate  
261 increases, as does the ionization efficiency, for many molecules.<sup>42</sup>

262

## 263 System 2

264 As mentioned earlier, System 2 aimed at approximating the behavior of the analyte on the surface of  
265 the droplet where most of the charged ions are expected to be located. While, the presence of  
266 hydronium is well established, we are confronted with the paradox that the most common ions formed  
267 in positive electrospray ionization are not  $[\text{M}+\text{H}_3\text{O}]^+$ , but rather  $[\text{M}+\text{H}]^+$ . This indicates that either  $\text{H}^+$   
268 ions also exist in the form of  $\text{H}^+$  alongside  $\text{H}_3\text{O}^+$  ions on the surface of the droplet or during the final

interactions of M and H<sub>3</sub>O<sup>+</sup> one of the protons disengages from the H<sub>2</sub>O molecule and forms the [M+H]<sup>+</sup> ion. In order to account for this discrepancy, in System 2, we describe the H<sup>+</sup> ions as freely floating ions that are not covalently bound to water molecules. In this case, H<sup>+</sup> was described in the same way as other ions like Na<sup>+</sup> and Cl<sup>-</sup> are described in GROMACS using the CGenFF force field. In this description, the mass of H<sup>+</sup> was set at 1.0080 g/mol and the charge (q) was set at +1. Since this is a single atom, the total charge (qtot) was also set at +1. As a point of reference, Na<sup>+</sup> ions in CGenFF are described as single atoms with a mass of 22.98977 g/mol, charge of +1 and a total charge of +1. A 4 x 4 x 4 nm solvent box was created with approximately 600 water molecules (TIP3P) and 600 H<sup>+</sup> ions. The number of 600 H<sup>+</sup> was determined based on pilot simulations so that H<sup>+</sup> would remain evenly distributed inside the box throughout the simulation to ensure continuous interactions with the analyte. A smaller number of ions resulted in the ions starting evenly distributed but during the simulation moving to the outer parts of the box and not sufficiently interacting with the chemical which often remained towards the center of the box (Figure S3-5 and Text S2).

282

### Simulation setup

The simulations were conducted using GROMACS version 2023.2. The simulation protocol started with a steepest descend minimization with 50,000 steps as the maximum number of minimization steps to perform and <1000 kJ/mol/nm as the threshold at which the minimization process can stop. The minimization and subsequent simulation steps were run using Verlet as the cut-off scheme for neighbor searching and Fast Smooth Particle-Mesh Ewald electrostatics (FSPME or PME in GROMACS) for modeling the electrostatic interactions. The short-range electrostatic cut-off points for Coulomb and van der Waals interactions were set to 1.2 nm which is recommended for CGenFF.<sup>43</sup> The temperature was set at 300 K and it was controlled with a Berendsen thermostat in NVT and a Parrinello-Rahman

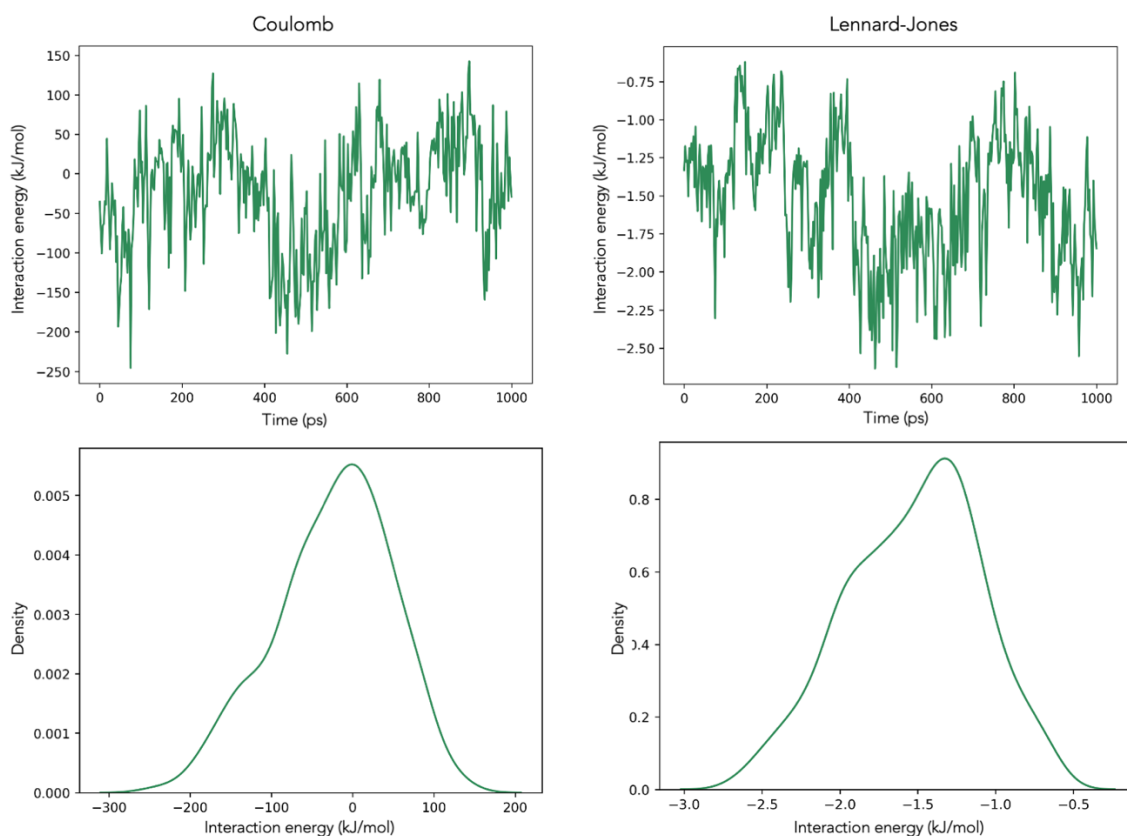
barostat in NPT. System equilibration was conducted in two stages, the NVT stage where volume and temperature were kept constant and the NPT stage where pressure and temperature were kept constant. The simulation step was set at 0.5 fs using a leapfrog integrator and the simulation length was 200 ps. The production step following equilibration was conducted using the same simulation step and integrator as previously but in this case the simulation length was 1000 ps (1 ns). All the mdp files for the minimization, equilibration and production steps with all the details are available on GitHub (<https://github.com/dimitriabrahamsson/electro-spray>).

299

## Calculation of interactions and model development

For both System 1 and System 2, we calculated the short-range Lennard-Jones and short-range Coulomb interactions between the analyte and each group of molecules in the system. For System 1, the sets were i) analyte and water, ii) analyte and methanol, and iii) analyte and  $\text{H}_3\text{O}^+$  ions. For System 2, the sets were i) analyte and water and ii) analyte and  $\text{H}^+$  ions, however, we only considered the set of analyte and  $\text{H}^+$  ions since the interactions with water were already described in System 1. As these are short-range interactions, it is important to point out that this includes only the molecules around the analyte that are within the short-range distance, which was set at 1.2 nm. The interactions were calculated using the gmx energy command in GROMACS. The generated files contained the interaction energies (kJ/mol) over time (ps) in the form of a time series. An example of the Coulomb and Lennard-Jones interactions for caffeine is shown in Figure 3. The top figures show the interactions over time and the bottom figures show the distribution of the observed interaction energies using the kernel density estimate.

313



315

316 Figure 3: Coulomb and Lennard-Jones interactions between caffeine and  $H^+$  ions during the simulation  
 317 of System 2. The top figures show the interactions over time and the bottom figures show the  
 318 distribution of the observed interaction energies.

319

320 Our model is based on the idea that the RRF of a given compound in ESI+ can be described as a  
 321 function of the Coulomb and Lennard-Jones interactions between the compound and all other  
 322 molecules in the solution. RRF was expressed as a function of the Coulomb and Lennard-Jones using a  
 323 multilinear regression model.

324

325 The model was as follows:

326



$$\log RRF = l LJ + c Coul + const \quad eq.6$$

where,

LJ is the Lennard-Jones interactions, Coul is the Coulomb interactions.

One critical challenge when incorporating these interactions into a model is finding which metrics are meaningful for the purposes of the model. We applied a Monte Carlo simulation approach to randomly sample the Coulomb and Lennard-Jones distributions (3 percentile points per distribution plus the standard deviation) 100 times for each set of molecules (i.e., System 1: i. analyte and water, ii. analyte and methanol, and iii. analyte and H<sub>3</sub>O<sup>+</sup> ions; System 2: analyte and H<sup>+</sup> ions) and selected the best performing model for each set.

Expanding the  $l LJ$  and  $c Coul$  terms of the equation we get:

$$l LJ = l_1 LJ_{p1} + l_2 LJ_{p2} + l_3 LJ_{p3} + l_4 LJ_{std} \quad eq.7$$

$$c Coul = c_1 Coul_{p1} + c_2 Coul_{p2} + c_3 Coul_{p3} + c_4 Coul_{std} \quad eq.8$$

where,

p1, p2 and p3 are the 3 percentile points from the distribution (e.g., 20, 50 and 70) and std is the standard deviation of the distribution.

The coefficients and the intercept of the model were determined through a least-squares minimization using the statsmodels package (version 0.14.0) in Python (version 3.10.11). The script is available on GitHub ( <https://github.com/dimitriabrahamsson/electro-spray> ). The model was evaluated by examining the R<sup>2</sup>, the mean absolute error (MAE) and the p-values of the coefficients. After selecting

the best performing model for each set of molecules from both systems, we built a composite model with the parameters whose p-values were lower than 0.1. We purposely chose a higher cutoff point at this stage in order to be more inclusive, however, in the final model, only the p-values below 0.05 were considered significant. The final model was evaluated based on its  $R^2$ , the mean absolute error (MAE) and the p-values of the coefficients. We also tested whether the addition of other physicochemical properties, i.e. vapor pressure ( $P_V$ ), water solubility ( $S_W$ ), the equilibrium partitioning ratio between octanol and water ( $K_{OW}$ ), air and water ( $K_{AW}$ ), methanol and water ( $K_{MW}$ ), methanol and air ( $K_{MA}$ ), and the innate charge of the molecule improved the predictive accuracy of the model.  $P_V$  and  $S_W$  were calculated with OPERA 2.6<sup>44</sup> available on the dashboard<sup>8</sup>.  $K_{OW}$ ,  $K_{AW}$ ,  $K_{MW}$  and  $K_{MA}$  were calculated with UFZ-LSER.<sup>45</sup> The innate charge of the molecule was determined by examining the structure and its protonation state at pH 5 and noting 0 if it was neutral, +1 (or more) if it had a positive innate charge, and -1 (or less) if it had a negative innate charge. Only the parameters whose coefficient showed a p-value of less than 0.05 were considered significant and were included in the model. A parameter with a mere increase in  $R^2$  without a significant p-value would not be included in the model.

The predictive power of the model was further evaluated with a 10-fold cross validation and a y-randomization. During the 10-fold cross-validation, the dataset was first divided into 10 equally sized sub-datasets. Then, during each fold one dataset was set as the testing set and the remaining sub-datasets were compiled into a training set. The model was trained on the training set and tested on the testing set. The process was repeated 10 times (10-fold). It is important to note at this point that when applying a k-fold cross validation and when dividing the dataset into training and testing there is always a possibility of encountering compounds in the testing set that are outside the applicability domain of the trained model. In order to account for this discrepancy, if a prediction was 2 log units higher than the highest value in the dataset or 2 log units lower than the lowest value in the dataset it was

373 considered outside the applicability domain and it was excluded from the evaluation. The compounds  
374 that were excluded from any particular fold of the cross-validation exercise were still included in the  
375 discussion section of the paper. The purpose of the k-fold cross validation is to evaluate the predictive  
376 power of the model outside its training set and to identify outlier compounds in the dataset. These  
377 compounds are considered outliers in the sense that they represent physicochemical properties that  
378 are dissimilar to the ones in the training set and in order for the model to make accurate predictions,  
379 they have to be included in the training set.

380 For the y-randomization, the y variable, in this case the RRF was randomly shuffled, and the model  
381 was developed as previously by dividing the dataset into training and testing sets. The process was  
382 repeated 5 times, and the predictions were averaged across the 5 iterations. The purpose of the y-  
383 randomization is to evaluate the extent to which the model predictions are different from random  
384 predictions. This helps to determine whether the model is making meaningful predictions and whether  
385 it has been overparametrized. The lower the  $R^2$  of the y-randomization and the more different it is from  
386 the  $R^2$  of the cross-validation, the higher the likelihood that the model is making meaningful predictions  
387 that are distinct from random predictions.

388 One of the challenges we encountered is that the generated CGenFF topologies often included high  
389 penalties (> 50) for a charge, a bond, an angle, a dihedral or an improper group. While high penalties  
390 do not necessarily mean large errors, they do denote a low similarity with the build-by-analogy  
391 structure in CGenFF and it is recommended to apply caution when using such structures because they  
392 may require further validation. This may be an important issue in the case of protein dynamics and  
393 ligand binding, however, in our case, it is unclear how these penalties or uncertainties may influence  
394 our calculations. To address this issue, we tested the robustness of the model by incrementally  
395 removing compounds with high penalties starting from the ones with the highest penalties to the ones

with the lowest penalties. This resulted in 10 different models with a different cutoff point as the maximum acceptable penalty ranging from 500 to 50. We then examined the changes in the  $R^2$  of the model as the number and type of chemicals in the dataset changed. In order to avoid introducing errors in the first steps of the model development, we developed the first iteration of the model with chemicals that had a penalty less than 300.

## RESULTS AND DISCUSSION

### Experimental measurements

The observed log RRF values for the chemicals in our dataset ranged from 1.73 to 3.17 with Cinchophen showing the lowest value and Thiabendazole showing the highest value (Supplementary Spreadsheet). As RRF is the ratio of abundance to concentration, higher RRF values indicate higher ionization efficiency (higher abundance at lower concentrations). This observation is in agreement with data from our previous study<sup>12</sup> where Cinchophen showed a lower RRF compared to Thiabendazole. It should be noted that the two studies use the same mixtures (in part) but different methods and different instruments (same type – Agilent LC-QTOF-MS – but different instrument). Despite the differences in methods and instrumentation, the differences in the RRF values of the two chemicals are preserved. This observation is supportive of the ionization efficiency (IE) scale approach developed by Oss et al.<sup>22</sup> where a set of RRF values can be represented as a scale of relative ionization efficiencies and that that scale should in principle be transferable across different methods.

### Model development

From the models developed for System 1, the best performing models for predicting log RRF showed an  $R^2$  of 0.42 when using the analyte-water interactions, 0.37 when using the analyte-methanol

interactions and 0.39 when using the analyte-H<sub>3</sub>O<sup>+</sup> interactions (Supplemental Spreadsheet). From the models developed for System 2, the best performing model for predicting log RRF showed an R<sup>2</sup> of 0.71 (Supplemental Spreadsheet). This observation indicates that the final stage of ionization [M+H]<sup>+</sup> is better predicted by the interactions of the analyte with the H<sup>+</sup> ions on the surface of the droplet (Figure 1 – step 2 of IEM) than by the interactions of the analyte with the other molecules while in the center of the droplet (Figure 1 – step 1 of IEM). This is not to say that there is no predictive value in the interactions of the analyte with the solvent molecules. Previous studies have demonstrated the impacts of different solvents on the ionization efficiency of small molecules<sup>46,47</sup> and this is in agreement with our calculations from System 1. This observation just indicates that the interactions of the analyte on the droplet are potentially more determining of the ionization efficiency of the analyte. The final composite model consisted of the following parameters. System 1: p2, p3 and the standard deviation for Lennard-Jones interactions between the analyte and water; System 2: p1, p2, p3 and the standard deviations for Lennard-Jones and Coulomb interactions between the analyte and H<sup>+</sup> ions. All coefficients for the abovementioned parameters showed p-values below 0.05 (Table S1). Out of all the physicochemical properties that we tested, the only one that showed a statistically significant contribution was the water solubility of the analyte (S<sub>w</sub>). The final model showed an R<sup>2</sup> of 0.82 and an MAE of 0.13.

442

443 The coefficients and intercept of the developed model were determined to be as follows:

444

445  $\log RRF = l LJ + c Coul - 0.14 S_W + 2.51$  *eq. 9*

446 where,

447

448  $l LJ = 0.63 LJ_{p1}^H - 5.80 LJ_{p2}^H + 4.92 LJ_{p3}^H + 3.49 LJ_{std}^H + 0.21 LJ_{p2}^W - 0.20 LJ_{p3}^W + 0.35 LJ_{std}^W$  *eq. 10*

449  $c Coul = -0.01 Coul_{p1}^H + 0.03 Coul_{p2}^H - 0.03 Coul_{p3}^H - 0.04 Coul_{std}^H$  *eq. 11*

450

451 where,

452  $LJ^H$  are the Lennard-Jones interactions between the analyte and  $H^+$  ions from System 2

453  $Coul^H$  are the Coulomb interactions between the analyte and  $H^+$  ions from System 2

454  $LJ^W$  are the Lennard-Jones interactions between the analyte and water molecules from System 1

455 the values for p1, p2 and p3 in System 2 were 0.5, 34 and 50

456 the values for p2, and p3 in System 1 were 44 and 89.

457

458 The p-values of the coefficients and the intercept were all below 0.05 (Table S1). The  $R^2$  and MAE of  
459 the model were comparable to those in the study of Oss et al.<sup>22</sup> where they observed an  $R^2$  of 0.67 and  
460 a standard residual error of 0.86 log units. While the two studies are very different in the  
461 computational approaches, they both share datasets of similar size (48 vs 62) and they both use  
462 multilinear regression models as their final predictive models and thus allowing for meaningful  
463 comparisons.

464 We examined whether the differences between the experimental and modeled RRF (absolute errors)  
465 could be explained due to the different retention times (RT) of the chemicals and by extension due to  
466 the different ratios of water to methanol, but we did not observe a significant association between the  
467 two. Neither did we observe a significant association between RRF and RT (Figure S6).

468 Our modeling calculations showed that the interactions of the analyte with the water molecules in  
469 System 1 were similar but slightly more predictive than the interactions of the analyte with the H<sub>3</sub>O<sup>+</sup>  
470 ions ( $R^2$ : 0.42 vs 0.39). Given the great collinearity of these two variables, including both of them in the  
471 model renders the coefficients for H<sub>3</sub>O<sup>+</sup> insignificant ( $p > 0.05$ ). This observation suggests that, at least  
472 in terms of interactions with the analyte, the H atoms in the H<sub>2</sub>O molecules are not distinguishable  
473 from the H atoms in the H<sub>3</sub>O<sup>+</sup> ions.

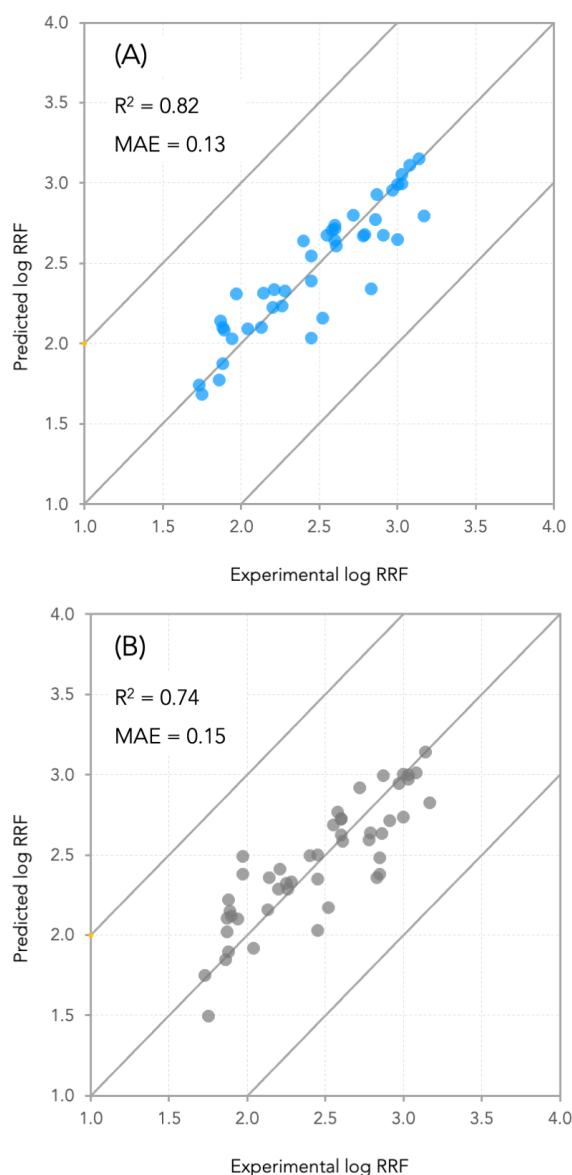
474 As mentioned earlier in the methods, we examined the effect that compounds with high penalties  
475 may have on the predictive power of the model. Including all compounds with penalties over 300  
476 resulted in a small decrease in  $R^2$  (0.82 vs 0.74) and a small increase in the MAE (0.13 vs 0.15) of the  
477 model (Figure 4). After incrementally removing compounds with high penalties from the dataset, we  
478 observed that the  $R^2$  of the model appeared to be consistent with an increase around cutoff points of  
479 250 and 300 (Figure S7), which confirmed our initial cut-off point of 300. We should note at this point  
480 that while in this particular case the effect of including compounds with high penalties appears to be  
481 minimal, we do not know how that may manifest in other datasets with different composition and with  
482 compounds with higher penalties.

483

484

485

486



488

489 Figure 4: Experimental and predicted values of log transformed RRF. The diagonal lines show the 1-to-1  
 490 agreement line, and the  $\pm 1$  log unit deviation line. Plot A shows the results of the model after  
 491 removing compounds with a penalty over 300 (dataset  $n = 42$ ). Six chemicals were excluded from the  
 492 dataset in this iteration. Plot B shows the results of the model including all the chemicals in the dataset  
 493 (dataset  $n = 48$ ).  $R^2$  is the coefficient of determination and MAE is the mean absolute error between the  
 494 predictions and the experimental values.



495

496 The 10-fold cross validation showed an  $R^2$  of 0.52 and an MAE of 0.25 for compounds that were not  
497 included in the training set (Figure S8A). This shows that the model can make reasonable predictions  
498 for chemicals that were not included in the training set. Two chemicals showed to be outside of the  
499 applicability domain of the model (based on the definition in the methods section). These two  
500 chemicals were Furalaxyl and Dicrotophos (Figure S9). During the cross-validation Furalaxyl showed an  
501 absolute error of 10.1 log units and Dicrotophos an absolute error of 5.58 log units. Both chemicals had  
502 penalties lower than 300 so it does not seem that their penalties would be a likely explanation  
503 (Supplemental Spreadsheet). Most likely is that these two chemicals are structurally and  
504 physicochemically distinct from the other chemicals in the dataset. This is supported by the observation  
505 that when these two compounds are included in the dataset their absolute errors are 0.001 log units  
506 for Furalaxyl and 0.03 log units for Dicrotophos (Figure S10).

507 The y-randomization showed that when the model is trained on random data the expected  $R^2$  is 0.03  
508 (Figure S8B). This is substantially lower than both the  $R^2$  of the model with all the chemicals (0.74) and  
509 the  $R^2$  of the 10-fold cross-validation (0.52). This observation suggests that the model is making  
510 meaningful predictions that are distinct from random predictions.

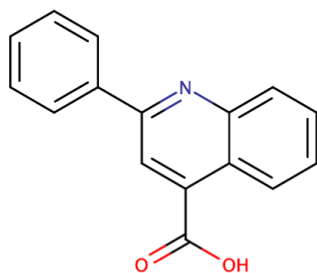
511 In trying to understand the contributions of the different interactions to RRF we examined the  
512 different terms of equation 9 for two chemicals that showed near 0 differences between experimental  
513 and predicted values of RRF. The two chemicals were 1) Cinchophen with a log RRF of 1.73, and 2)  
514 Loratadine with a log RRF of 3.15. Both chemicals' log RRF is determined to a larger extent by the  
515 Lennard-Jones and Coulomb interactions and to a smaller extent by their water solubility. For both  
516 chemicals, the Lennard-Jones interactions appear to have a positive contribution to RRF while the  
517 Coulomb interactions appear to have a negative contribution (Figure 5). This is consistent for all

chemicals in the dataset. When comparing Cinchophen and Loratadine, it appears that the lower RRF of Cinchophen is due to smaller  $l$  LJ and  $c$  Coul terms (Figure 5). The Lennard-Jones potential approximates the van der Waals interactions and the Coulomb potential represents the ability to engage in hydrogen bonding. Previous studies have suggested that increased non-polar character, which would be represented by the Lennard-Jones potential are associated with higher RRF, while increased polar character which would be represented by the Coulomb potential is associated with a decrease in RRF.<sup>22,23,27,48,49</sup> Our observations appear to be in agreement with the findings from previous studies.

Water solubility appears to play a small (Figure 5) yet significant role (Tables S1 and S2) in the model. For all compounds in the dataset, the  $s$   $S_w$  term has a positive contribution to RRF. Based on this observation, one would expect that compounds with higher water solubility would have a higher RRF. However, given that the term  $s$   $S_w$  is several orders of magnitude smaller than the  $l$  LJ and  $c$  Coul terms, the influence of  $s$   $S_w$  on RRF is minimal in comparison. In our developed model,  $s$   $S_w$  acts as a corrective factor rather than a determining factor. Furthermore, water solubility is known to decrease with increasing molecular weight<sup>50</sup>, which is also what we observed in our dataset. The contribution of the  $s$   $S_w$  for Cinchophen is slightly larger than that of Loratadine which is in agreement with their molecular weights (Cinchophen: 249.26 g/mol, Loratadine: 382.89 g/mol).

534

Cinchophen



$$\log RRF = 0.88 - 1.66 + 8.7 \times 10^{-5} + 2.51 = 1.74$$

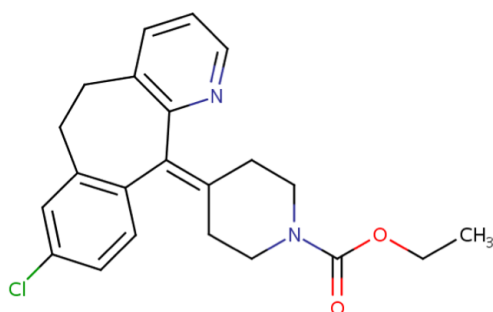
|  
*l LJ*

|  
*c Coul*

|  
*s S<sub>W</sub>*

|  
*const*

Loratadine



$$\log RRF = 1.39 - 0.75 + 1.1 \times 10^{-6} + 2.51 = 3.15$$

|  
*l LJ*

|  
*c Coul*

|  
*s S<sub>W</sub>*

|  
*const*

Figure 5: Contributions of Lennard-Jones interactions, Coulomb interactions and water solubility to the calculated log RRF for two compounds that showed near 0 errors between the experimental values and the predictions of log RRF.

## Limitations and future considerations

One limitation that needs to be acknowledged is that while our model showed good accuracy ( $R^2$ : 0.82) the computational cost of our approach is much higher than other approaches that rely on simpler descriptor generators like PaDEL<sup>51</sup> or Mordred<sup>52</sup>. Running one simulation on one system for

one chemical takes approximately 13 min using a ASUS GeForce GTX 1080 TI 11GB Turbo GPU. Conducting simulations for 48 chemicals, 2 systems and 3 replicates each comes up to 62.4h. This may limit the ability of the model to be used as an online application as it would require access to GPUs. The workflow is, however, applicable in PCs with NVIDIA GPUs.

Another limitation that should be acknowledged is that, in this study, we examined only one type of force field (CGenFF). Future applications could examine whether using other types of force fields like gaff2 from AMBER and GROMOS from GROMACS can produce better predictions than CGenFF.

Finally, on the experimental side, it should be acknowledged that for the purposes of this study, we tested only two solvents for our LC gradient, HPLC water and methanol. As the ionization efficiency of chemicals is known to vary by different solvents<sup>53</sup>, the effect of that variability on the modeling calculations is something that needs to be investigated further.

## CONCLUSION

Our study presents a novel approach for modeling the ionization efficiency of organic molecules. Our approach can be used in combination with existing approaches for concentration estimates of chemical compounds in environmental and biological samples. While there is a variety of modeling approaches for RRF, our view is that these approaches are complementary rather than competing. When trying to estimate concentrations of chemicals in environmental or biological samples, combining the results of multiple different approaches can help establish multiple layers of evidence that can be used in support of a prediction when analytical standards are unavailable.

565 **SUPPORTING INFORMATION**

566 The manuscript is accompanied by a Supporting Information pdf file and a Supplemental Spreadsheet  
567 containing the underlying data used in the paper.

568

569 **DECLARATIONS**

570 AVAILABILITY OF DATA AND MATERIALS

571 All code, underlying datasets, and chemical structures used in this study are publicly available and  
572 can be accessed on GitHub under the following repository:

573 <https://github.com/dimitriabrahamsson/electro-spray>

574 All python scripts are accompanied by instructions on how to run them in order to replicate the  
575 findings of the study.

576

577 **COMPETING INTERESTS**

578 The authors declare that they have no competing interests to report.

579

580 **FUNDING**

581 This study was funded by NIH/NIEHS (Grant Nos. R00ES032892, K99ES032892).

582

583 **AUTHOR CONTRIBUTIONS**

584 D.A. co-wrote the manuscript, designed the study, conducted the lab experiments, ran part of the  
585 simulations and worked on the modeling section. L-A.K co-wrote the manuscript, ran simulations and  
586 assisted with the modeling section. T.J. co-wrote the manuscript and helped with proof-reading and  
587 providing expertise on the analytical side. D.L. co-wrote the manuscript, supervised the experimental

part of the study, provided advice and expertise on designing and executing the experiments, and on interpreting the findings.

## ACKNOWLEDGEMENTS

We would like to thank Elin Ulrich from the U.S. EPA for assisting with providing the chemical mixtures that were used in this study. We would also like to thank June-Soo Park of the California EPA for providing lab space and equipment to support our experiments.

## REFERENCES

- (1) Newton, S. R.; McMahan, R. L.; Sobus, J. R.; Mansouri, K.; Williams, A. J.; McEachran, A. D.; Strynar, M. J. Suspect Screening and Non-Targeted Analysis of Drinking Water Using Point-of-Use Filters. *Environ. Pollut.* **2018**, *234*, 297–306. <https://doi.org/10.1016/j.envpol.2017.11.033>.
- (2) Wang, A.; Abrahamsson, D. P.; Jiang, T.; Wang, M.; Morello-Frosch, R.; Park, J.-S.; Sirota, M.; Woodruff, T. J. Suspect Screening, Prioritization, and Confirmation of Environmental Chemicals in Maternal-Newborn Pairs from San Francisco. *Environ. Sci. Technol.* **2021**, *55* (8), 5037–5049. <https://doi.org/10.1021/acs.est.0c05984>.
- (3) Panagopoulos Abrahamsson, D.; Wang, A.; Jiang, T.; Wang, M.; Siddharth, A.; Morello-Frosch, R.; Park, J.-S.; Sirota, M.; Woodruff, T. J. A Comprehensive Non-Targeted Analysis Study of the Prenatal Exposome. *Environ. Sci. Technol.* **2021**, *55* (15), 10542–10557. <https://doi.org/10.1021/acs.est.1c01010>.
- (4) Zhu, Y.; Barupal, D. K.; Ngo, A. L.; Quesenberry, C. P.; Feng, J.; Fiehn, O.; Ferrara, A. Predictive Metabolomic Markers in Early to Mid-Pregnancy for Gestational Diabetes Mellitus: A Prospective Test and Validation Study. *Diabetes* **2022**, *71* (8), 1807–1817. <https://doi.org/10.2337/db21-1093>.
- (5) Moschet, C.; Anumol, T.; Lew, B. M.; Bennett, D. H.; Young, T. M. Household Dust as a Repository of Chemical Accumulation: New Insights from a Comprehensive High-Resolution Mass Spectrometric Study. *Environ. Sci. Technol.* **2018**, *52* (5), 2878–2887. <https://doi.org/10.1021/acs.est.7b05767>.
- (6) Nuñez, J. R.; Colby, S. M.; Thomas, D. G.; Tfaily, M. M.; Tolic, N.; Ulrich, E. M.; Sobus, J. R.; Metz, T. O.; Teeguarden, J. G.; Renslow, R. S. Evaluation of In Silico Multifeature Libraries for Providing Evidence for the Presence of Small Molecules in Synthetic Blinded Samples. *J. Chem. Inf. Model.* **2019**, *59* (9), 4052–4060. <https://doi.org/10.1021/acs.jcim.9b00444>.
- (7) *Legal obstacles to toxic chemical research | Science*. <https://www.science.org/doi/10.1126/science.abl4383> (accessed 2023-11-30).
- (8) *CompTox Chemicals Dashboard*. <https://comptox.epa.gov/dashboard/> (accessed 2023-11-30).
- (9) Abrahamsson, D.; Siddharth, A.; Young, T. M.; Sirota, M.; Park, J.-S.; Martin, J. W.; Woodruff, T. J. In Silico Structure Predictions for Non-Targeted Analysis: From Physicochemical Properties to Molecular Structures. *J. Am. Soc. Mass Spectrom.* **2022**, *33* (7), 1134–1147. <https://doi.org/10.1021/jasms.1c00386>.

- (10) Abrahamsson, D.; Brueck, C. L.; Prasse, C.; Lambropoulou, D. A.; Koronaiou, L.-A.; Wang, M.; Park, J.-S.; Woodruff, T. J. Extracting Structural Information from Physicochemical Property Measurements Using Machine Learning—A New Approach for Structure Elucidation in Non-Targeted Analysis. *Environ. Sci. Technol.* **2023**, *57* (40), 14827–14838. <https://doi.org/10.1021/acs.est.3c03003>.
- (11) Liigand, P.; Liigand, J.; Kaupmees, K.; Kruve, A. 30 Years of Research on ESI/MS Response: Trends, Contradictions and Applications. *Anal. Chim. Acta* **2021**, *1152*, 238117. <https://doi.org/10.1016/j.aca.2020.11.049>.
- (12) Panagopoulos Abrahamsson, D.; Park, J.-S.; Singh, R. R.; Sirota, M.; Woodruff, T. J. Applications of Machine Learning to In Silico Quantification of Chemicals without Analytical Standards. *J. Chem. Inf. Model.* **2020**, *60* (6), 2718–2727. <https://doi.org/10.1021/acs.jcim.9b01096>.
- (13) Groff, L. C.; Grossman, J. N.; Kruve, A.; Minucci, J. M.; Lowe, C. N.; McCord, J. P.; Kapraun, D. F.; Phillips, K. A.; Purucker, S. T.; Chao, A.; Ring, C. L.; Williams, A. J.; Sobus, J. R. Uncertainty Estimation Strategies for Quantitative Non-Targeted Analysis. *Anal. Bioanal. Chem.* **2022**, *414* (17), 4919–4933. <https://doi.org/10.1007/s00216-022-04118-z>.
- (14) Johnson, T. A.; Abrahamsson, D. P. Quantification of Chemicals in Non-Targeted Analysis without Analytical Standards – Understanding the Mechanism of Electrospray Ionization and Making Predictions. *Curr. Opin. Environ. Sci. Health* **2023**, 100529. <https://doi.org/10.1016/j.coesh.2023.100529>.
- (15) Calixte, E. I.; Liyanage, O. T.; Kim, H. J.; Ziperman, E. D.; Pearson, A. J.; Gallagher, E. S. Release of Carbohydrate–Metal Adducts from Electrospray Droplets: Insight into Glycan Ionization by Electrospray. *J. Phys. Chem. B* **2020**, *124* (3), 479–486. <https://doi.org/10.1021/acs.jpccb.9b10369>.
- (16) Mandra, V. J.; Kouskoura, M. G.; Markopoulou, C. K. Using the Partial Least Squares Method to Model the Electrospray Ionization Response Produced by Small Pharmaceutical Molecules in Positive Mode. *Rapid Commun. Mass Spectrom.* **2015**, *29* (18), 1661–1675. <https://doi.org/10.1002/rcm.7263>.
- (17) Golubović, J.; Birkemeyer, C.; Protić, A.; Otašević, B.; Zečević, M. Structure–Response Relationship in Electrospray Ionization–Mass Spectrometry of Sartans by Artificial Neural Networks. *J. Chromatogr. A* **2016**, *1438*, 123–132. <https://doi.org/10.1016/j.chroma.2016.02.021>.
- (18) Henriksen, T.; Juhler, R. K.; Svensmark, B.; Cech, N. B. The Relative Influences of Acidity and Polarity on Responsiveness of Small Organic Molecules to Analysis with Negative Ion Electrospray Ionization Mass Spectrometry (ESI-MS). *J. Am. Soc. Mass Spectrom.* **2005**, *16* (4), 446–455. <https://doi.org/10.1016/j.jasms.2004.11.021>.
- (19) Kruve, A.; Kaupmees, K.; Liigand, J.; Leito, I. Negative Electrospray Ionization via Deprotonation: Predicting the Ionization Efficiency. *Anal. Chem.* **2014**, *86* (10), 4822–4830. <https://doi.org/10.1021/ac404066v>.
- (20) Liigand, P.; Kaupmees, K.; Haav, K.; Liigand, J.; Leito, I.; Girod, M.; Antoine, R.; Kruve, A. Think Negative: Finding the Best Electrospray Ionization/MS Mode for Your Analyte. *Anal. Chem.* **2017**, *89* (11), 5665–5668. <https://doi.org/10.1021/acs.analchem.7b00096>.
- (21) Alymatiri, C. M.; Kouskoura, M. G.; Markopoulou, C. K. Decoding the Signal Response of Steroids in Electrospray Ionization Mode (ESI-MS). *Anal. Methods* **2015**, *7* (24), 10433–10444. <https://doi.org/10.1039/C5AY02839F>.
- (22) Oss, M.; Kruve, A.; Herodes, K.; Leito, I. Electrospray Ionization Efficiency Scale of Organic Compounds. *Anal. Chem.* **2010**, *82* (7), 2865–2872. <https://doi.org/10.1021/ac902856t>.
- (23) Kruve, A.; Kaupmees, K. Predicting ESI/MS Signal Change for Anions in Different Solvents. *Anal. Chem.* **2017**, *89* (9), 5079–5086. <https://doi.org/10.1021/acs.analchem.7b00595>.



- (24) Nguyen, T. B.; Nizkorodov, S. A.; Laskin, A.; Laskin, J. An Approach toward Quantification of Organic Compounds in Complex Environmental Samples Using High-Resolution Electrospray Ionization Mass Spectrometry. *Anal. Methods* **2012**, 5 (1), 72–80. <https://doi.org/10.1039/C2AY25682G>.
- (25) Wu, L.; Wu, Y.; Shen, H.; Gong, P.; Cao, L.; Wang, G.; Hao, H. Quantitative Structure–Ion Intensity Relationship Strategy to the Prediction of Absolute Levels without Authentic Standards. *Anal. Chim. Acta* **2013**, 794, 67–75. <https://doi.org/10.1016/j.aca.2013.07.034>.
- (26) Liigand, P.; Liigand, J.; Cuyckens, F.; Vreeken, R. J.; Kruve, A. Ionisation Efficiencies Can Be Predicted in Complicated Biological Matrices: A Proof of Concept. *Anal. Chim. Acta* **2018**, 1032, 68–74. <https://doi.org/10.1016/j.aca.2018.05.072>.
- (27) Bieber, S.; Letzel, T.; Kruve, A. Electrospray Ionization Efficiency Predictions and Analytical Standard Free Quantification for SFC/ESI/HRMS. *J. Am. Soc. Mass Spectrom.* **2023**, 34 (7), 1511–1518. <https://doi.org/10.1021/jasms.3c00156>.
- (28) Palm, E.; Kruve, A. Machine Learning for Absolute Quantification of Unidentified Compounds in Non-Targeted LC/HRMS. *Molecules* **2022**, 27 (3), 1013. <https://doi.org/10.3390/molecules27031013>.
- (29) Liigand, J.; Wang, T.; Kellogg, J.; Smedsgaard, J.; Cech, N.; Kruve, A. Quantification for Non-Targeted LC/MS Screening without Standard Substances. *Sci. Rep.* **2020**, 10 (1), 5808. <https://doi.org/10.1038/s41598-020-62573-z>.
- (30) Daub, C. D.; Cann, N. M. How Are Completely Desolvated Ions Produced in Electrospray Ionization: Insights from Molecular Dynamics Simulations. *Anal. Chem.* **2011**, 83 (22), 8372–8376. <https://doi.org/10.1021/ac202103p>.
- (31) Kim, D.; Wagner, N.; Wooding, K.; Clemmer, D. E.; Russell, D. H. Ions from Solution to the Gas Phase: A Molecular Dynamics Simulation of the Structural Evolution of Substance P during Desolvation of Charged Nanodroplets Generated by Electrospray Ionization. *J. Am. Chem. Soc.* **2017**, 139 (8), 2981–2988. <https://doi.org/10.1021/jacs.6b10731>.
- (32) Luan, M.; Hou, Z.; Zhang, B.; Ma, L.; Yuan, S.; Liu, Y.; Huang, G. Inter-Domain Repulsion of Dumbbell-Shaped Calmodulin during Electrospray Ionization Revealed by Molecular Dynamics Simulations. *Anal. Chem.* **2023**, 95 (23), 8798–8806. <https://doi.org/10.1021/acs.analchem.2c05630>.
- (33) Luan, M.; Hou, Z.; Huang, G. Suppression of Protein Structural Perturbations in Native Electrospray Ionization during the Final Evaporation Stages Revealed by Molecular Dynamics Simulations. *J. Phys. Chem. B* **2022**, 126 (1), 144–150. <https://doi.org/10.1021/acs.jpccb.1c09130>.
- (34) Konermann, L.; Metwally, H.; McAllister, R. G.; Popa, V. How to Run Molecular Dynamics Simulations on Electrospray Droplets and Gas Phase Proteins: Basic Guidelines and Selected Applications. *Methods* **2018**, 144, 104–112. <https://doi.org/10.1016/j.ymeth.2018.04.010>.
- (35) Ulrich, E. M.; Sobus, J. R.; Grulke, C. M.; Richard, A. M.; Newton, S. R.; Strynar, M. J.; Mansouri, K.; Williams, A. J. EPA’s Non-Targeted Analysis Collaborative Trial (ENTACT): Genesis, Design, and Initial Findings. *Anal. Bioanal. Chem.* **2019**, 411 (4), 853–866. <https://doi.org/10.1007/s00216-018-1435-6>.
- (36) Chemaxon. <https://chemaxon.com/> (accessed 2024-06-01).
- (37) Van Berkel, G. J.; Zhou, F.; Aronson, J. T. Changes in Bulk Solution pH Caused by the Inherent Controlled-Current Electrolytic Process of an Electrospray Ion Source. *Int. J. Mass Spectrom. Ion Process.* **1997**, 162 (1), 55–67. [https://doi.org/10.1016/S0168-1176\(96\)04476-X](https://doi.org/10.1016/S0168-1176(96)04476-X).
- (38) CHARMM-GUI. <https://www.charmm-gui.org/> (accessed 2024-01-21).



- (39) *Periodic boundary conditions - GROMACS 2024.2 documentation.*  
<https://manual.gromacs.org/current/reference-manual/algorithms/periodic-boundary-conditions.html> (accessed 2024-06-01).
- (40) Wolf, M. G.; Groenhof, G. Explicit Proton Transfer in Classical Molecular Dynamics Simulations. *J. Comput. Chem.* **2014**, *35* (8), 657–671. <https://doi.org/10.1002/jcc.23536>.
- (41) Smith, J. N.; Flagan, R. C.; Beauchamp, J. L. Droplet Evaporation and Discharge Dynamics in Electrospray Ionization. *J. Phys. Chem. A* **2002**, *106* (42), 9957–9967.  
<https://doi.org/10.1021/jp025723e>.
- (42) Calixte, E. I.; Liyanage, O. T.; Gass, D. T.; Gallagher, E. S. Formation of Carbohydrate–Metal Adducts from Solvent Mixtures during Electrospray: A Molecular Dynamics and ESI-MS Study. *J. Am. Soc. Mass Spectrom.* **2021**, *32* (12), 2738–2745. <https://doi.org/10.1021/jasms.1c00179>.
- (43) Vanommeslaeghe, K.; Hatcher, E.; Acharya, C.; Kundu, S.; Zhong, S.; Shim, J.; Darian, E.; Guvench, O.; Lopes, P.; Vorobyov, I.; Mackerell Jr., A. D. CHARMM General Force Field: A Force Field for Drug-like Molecules Compatible with the CHARMM All-Atom Additive Biological Force Fields. *J. Comput. Chem.* **2010**, *31* (4), 671–690. <https://doi.org/10.1002/jcc.21367>.
- (44) Mansouri, K.; Grulke, C. M.; Judson, R. S.; Williams, A. J. OPERA Models for Predicting Physicochemical Properties and Environmental Fate Endpoints. *J. Cheminformatics* **2018**, *10* (1), 10. <https://doi.org/10.1186/s13321-018-0263-1>.
- (45) *UFZ - LSER Database.*  
[https://www.ufz.de/index.php?en=31698&contentonly=1&m=0&lserd\\_data\[mvc\]=Public/start](https://www.ufz.de/index.php?en=31698&contentonly=1&m=0&lserd_data[mvc]=Public/start) (accessed 2024-05-26).
- (46) Krueve, A. Influence of Mobile Phase, Source Parameters and Source Type on Electrospray Ionization Efficiency in Negative Ion Mode. *J. Mass Spectrom.* **2016**, *51* (8), 596–601. <https://doi.org/10.1002/jms.3790>.
- (47) Liigand, J.; Krueve, A.; Leito, I.; Girod, M.; Antoine, R. Effect of Mobile Phase on Electrospray Ionization Efficiency. *J. Am. Soc. Mass Spectrom.* **2014**, *25* (11), 1853–1861. <https://doi.org/10.1007/s13361-014-0969-x>.
- (48) Cech, N. B.; Enke, C. G. Relating Electrospray Ionization Response to Nonpolar Character of Small Peptides. *Anal. Chem.* **2000**, *72* (13), 2717–2723. <https://doi.org/10.1021/ac9914869>.
- (49) Hermans, J.; Ongay, S.; Markov, V.; Bischoff, R. Physicochemical Parameters Affecting the Electrospray Ionization Efficiency of Amino Acids after Acylation. *Anal. Chem.* **2017**, *89* (17), 9159–9166. <https://doi.org/10.1021/acs.analchem.7b01899>.
- (50) Tolls, J.; van Dijk, J.; Verbruggen, E. J. M.; Hermens, J. L. M.; Loeprecht, B.; Schüürmann, G. Aqueous Solubility–Molecular Size Relationships: A Mechanistic Case Study Using C10- to C19-Alkanes. *J. Phys. Chem. A* **2002**, *106* (11), 2760–2765. <https://doi.org/10.1021/jp011755a>.
- (51) *PaDEL-descriptor: An open source software to calculate molecular descriptors and fingerprints - Yap - 2011 - Journal of Computational Chemistry - Wiley Online Library.*  
<https://onlinelibrary.wiley.com/doi/full/10.1002/jcc.21707> (accessed 2024-06-01).
- (52) Moriwaki, H.; Tian, Y.-S.; Kawashita, N.; Takagi, T. Mordred: A Molecular Descriptor Calculator. *J. Cheminformatics* **2018**, *10* (1), 4. <https://doi.org/10.1186/s13321-018-0258-y>.