

1 **Critical Assessment of pH-Dependent Lipophilic Profiles of Small Molecules: Which**
2 **One Should We Use and In Which Cases?**

3 Esteban Bertsch¹⁺, Sebastian Suñer¹⁺, Silvana Pinheiro^{1,2}, William J. Zamora^{1,2,3,*}

- 4 1. CBio³ Laboratory, School of Chemistry, University of Costa Rica, San Pedro, San José, Costa
5 Rica
6 2. Laboratory of Computational Toxicology and Artificial Intelligence (LaToxCIA), Biological
7 Testing Laboratory (LEBi), University of Costa Rica, San Pedro, San José, Costa Rica.
8 3. Advanced Computing Lab (CNCA), National High Technology Center (CeNAT), Pavas, San
9 13 José, Costa Rica

10

11 * Corresponding author: william.zamoraramirez@ucr.ac.cr (WJZ)

12 ⁺These authors contributed equally

13

14 **Keywords** Partition coefficient, Lipophilic profiles, Machine learning, Chemoinformatics, Drug
15 Design, pH.

16

17

18

19

20

21

22

23

24

25

26

27

28 Abstract

29 Lipophilicity is a physicochemical property with wide relevance in drug design and also
30 applied in areas such as food chemistry, environmental chemistry, and computational biology. This
31 descriptor strongly influences the absorption, distribution, permeability, bioaccumulation, protein-
32 binding, and biological activity of bioorganic compounds. Lipophilicity is commonly expressed
33 as the *n*-octanol/water partition coefficient (P_N) for neutral molecules, whereas for molecules with
34 ionizable groups, the distribution coefficient (D) at a given pH is used. The $\log D_{\text{pH}}$ is usually
35 predicted using a pH correction over the $\log P_N$ using the $\text{p}K_a$ of ionizable molecules, while often
36 ignoring the apparent ionic partition (P_I^{app}) because of the challenge of predicting the partitioning
37 of the charged species and/or related species (e.g., ion-pairs, counterions, molecular aggregates).
38 In this work, we studied the impact of P_I^{app} on the prediction of lipophilicity of small molecules by
39 modeling 225 $\log D_{\text{pH}}$ of a set of experimental values using the formalism that takes into account a
40 pH correction (see Eq. 1) and the one considering the apparent partition of ionic species (see Eq.
41 2). Our findings show that a better fit is obtained by considering the apparent ionic partition while
42 ignoring its contribution can lead to inadequate computational predictions. In this context, we
43 developed machine learning algorithms to determine in which cases the P_I^{app} should be considered.
44 The results indicate that small, rigid, and unsaturated molecules with $\log P_N$ close to zero which
45 present a significant proportion of ionic species in the aqueous phase, were better modeled using
46 Eq. 2. In addition, we validated our findings using a test and two external set which include small
47 molecules and amino acids analogs where the logistic regressions, random forest classifications,
48 and support vector machine models predicted the better formalism to determine the $\log D_{\text{pH}}$ for
49 each molecule with high accuracies, sensitivities, and specificities. Finally, our findings can serve
50 as guidance to the scientific community working in early-stage drug design, food, and
51 environmental chemistry who deal with ionizable molecules, to determine a priori which pH-
52 dependent lipophilicity profile should be used depending on the structure of a substance in their
53 research.

54

55

56

57 Introduction

58 Lipophilicity has been a relevant physicochemical property in pharmaceutical research
59 since the late 1800s, where the toxicity and anesthetic properties of several substances have been
60 correlated to their solubilities in water and oil/water partition coefficients.¹ In addition, this
61 property has been associated with several pharmacokinetic properties, such as enzyme binding²,
62 toxicity³, solubility⁴, membrane permeability⁵, and bioaccumulation.⁶ Thus, lipophilicity has been
63 considered a significant descriptor in drug discovery metrics, such as Lipinski's⁷ and Veber's⁸
64 empirical rules, which are intended to optimize oral bioavailability for drug-like compounds. The
65 partition coefficient (P_N) describes the equilibrium of a molecule between the organic and aqueous
66 phases, where the *n*-octanol/water system has historically been the medium of choice in
67 pharmaceutical research because of its high correlation with biological activities.^{9,10} However,
68 $\log P_N$ only describes the equilibrium of molecules in their neutral states, which implies an
69 unrealistic protonation state for most molecules with ionizable groups at physiological pH.

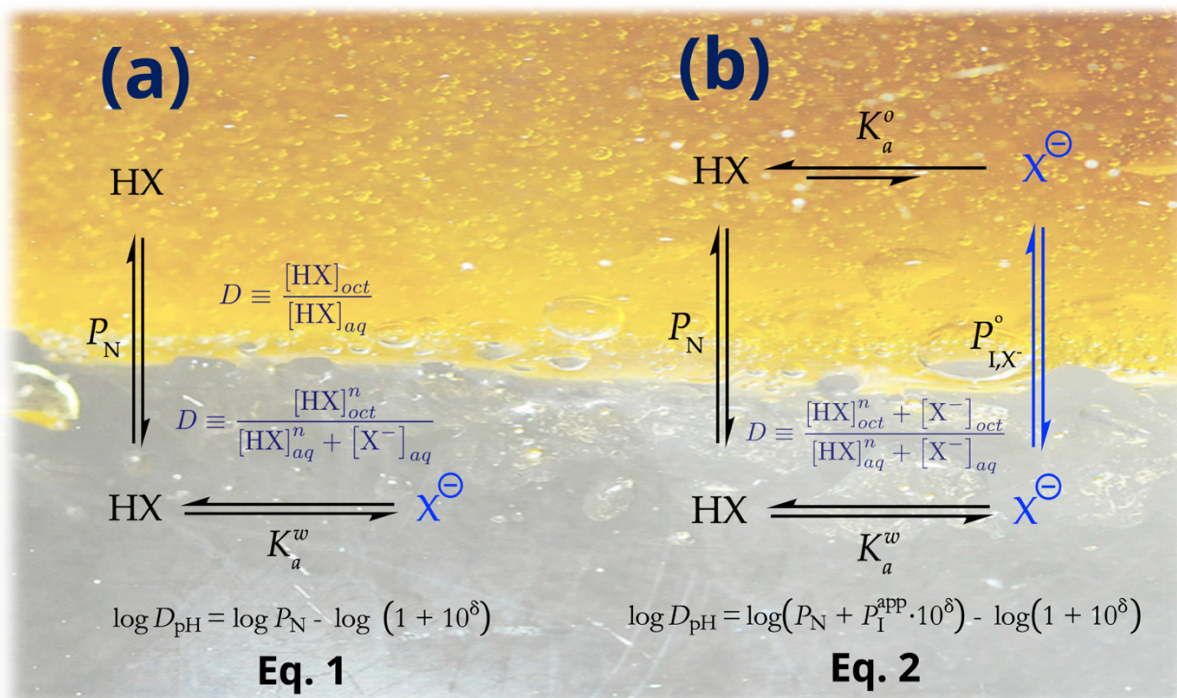
70 Since the pH of the solution directly affects the concentration of neutral and ionic species,
71 the equilibrium constant varies with pH, which also means that the lipophilicity of a compound is
72 dependent on it. The partition coefficient as a function of pH is often called distribution coefficient
73 ($\log D_{\text{pH}}$).¹¹ The $\log D_{\text{pH}}$ is often a more proper descriptor for human bioavailability due to the
74 frequent pH-dependence of drugs. This property has shown be useful in QSAR models to explain
75 how small molecules have human brain cells permeability¹² or binding to human serum albumin¹³.
76 The $\log D_{\text{pH}}$ has also been used as an effective predictor of pH-dependent lipophilicity profiles for
77 small molecules¹⁴ and to characterize structural properties in proteins and peptides, such as
78 protein-folding and aggregation¹⁵, solubility¹⁶, and antimicrobial activity^{17,18}, through pH-
79 dependent lipophilicity scales.^{19,20}

80 As an alternative to the experimentally determined $\log D_{\text{pH}}$ values, theoretical lipophilicity
81 profiles give the opportunity to obtain this descriptor quickly and often with high accuracy.^{14,21,22}
82 Equation 1 models $\log D_{\text{pH}}$ as a function of pH for monoacidic and monobasic compounds. This
83 equation is derived as a mass balance between ionic and neutral species in thermodynamic
84 equilibrium in the aqueous phase. This model assumes that the organic phase holds mostly neutral
85 species, so that the acid-base dissociation is negligible, and it also assumes that there is not a
86 partition equilibrium for the ionic species.²³

87
$$\log D_{\text{pH}} = \log P_{\text{N}} - \log (1 + 10^{\delta}) \quad [1]$$

88 where $\delta = \text{pH} - \text{p}K_{\text{a}}$ for acids, and $\delta = \text{p}K_{\text{a}} - \text{pH}$ for bases.

89
 90 Figure 1a displays the equilibria from which it is derived. Eq. 1 has been used to easily
 91 calculate $\log D_{\text{pH}}$ from $\log P_{\text{N}}$ values obtained by empirical computational models.²⁴⁻²⁶ This equation
 92 was widely used in $\log D_{\text{pH}}$ estimation methods in the SAMPL6 and SAMPL7 blind challenge,
 93 which is a large-scale comparative evaluation for drug design predictive models.^{27,28}



94
 95 **Figure 1.** Representations of the mechanism of partition for a symbolic ionizable acidic molecule
 96 for both neutral (HX) and ionic (X⁻) species using (a) Equation 1 and (b) Equation 2. The
 97 theoretical partition of the charged organic specie (P_{I,X^-}) has been replaced by experimental
 98 measurable apparent partitioning ($P_{\text{I}}^{\text{app}}$) in Eq. 2.

99
 100 Equation 2 represents the extended lipophilicity profile for monoprotic acids and bases (see
 101 Figure 1b). This model considers acid-base ionization in both water and *n*-octanol phases where
 102 ionic species migrate between the phases.

103
$$\log D_{\text{pH}} = \log(P_{\text{N}} + P_{\text{I}}^{\text{app}} \cdot 10^{\delta}) - \log(1 + 10^{\delta})$$
 [2]

104

105 Equation 2 is commonly called ionic partition P_{I} model²⁹, which represents a simplification
106 that only considers the partition of the charged organic specie (see Figure 1b). Experimental
107 techniques for lipophilicity evaluation such as shake-flask, potentiometric, and chromatographic
108 methods³⁰, can measure but do not allow direct identification of the nature of the ionic specie (es)
109 involved in the partitioning, hence, the partition of ionic species is measured as an apparent
110 partitioning ($P_{\text{I}}^{\text{app}}$). This experimentally measurable apparent partition coefficient depends on the
111 background salt³¹, compound concentration³², and may involve much more complex species such
112 as ion-pairs³³⁻⁴⁰, and aggregates⁴¹. Some studies have simplified the $P_{\text{I}}^{\text{app}}$ to the partition of only
113 ionic organic species (P_{I}) because these methods used have been parametrized by using
114 experimental $P_{\text{I}}^{\text{app}}$ values^{14,42}, while other theoretical studies have modeled it using the
115 participation of ion-pairs (P_{IP})^{21,22}. Recently, an alternative model¹⁴ to that of ion-pair partitioning
116 has been used by applying the theory of ionic transfer between two immiscible electrolyte
117 solutions (ITIES)^{43,44}, obtaining excellent predictions of experimental $\log D_{\text{pH}}$ values. Previous
118 experimental trials have also shown the importance of the $P_{\text{I}}^{\text{app}}$ of ionizable molecules in *n*-
119 octanol/water systems³³⁻⁴⁰. Recently, Disdier *et al.* measured the $\log D_{\text{pH}}$ at different pH values of
120 a set of 13 compounds via the shake-flask method⁴⁵, where they fitted their experimental values to
121 lipophilicity formalisms for mono- and poly substituted acids, amphoteric, and zwitterionic species
122 derived on previous theoretical studies.⁴⁶ The relevance of $P_{\text{I}}^{\text{app}}$ for small ionic molecules between
123 aqueous and organic phases has also been studied through interphase transfer mechanisms of
124 substances via ionic partition diagrams as a function of pH obtained through cyclic voltammetry.⁴⁷⁻
125 ⁴⁹

126 Despite the lack of a consensus formalism to model $\log D_{\text{pH}}$ as a function of the $P_{\text{I}}^{\text{app}}$ and
127 considering that different theoretical approaches have shown similar trends^{14,21,22}, Equation 2 has
128 been successfully used for modeling lipophilicity of ionized compounds in many areas of basic
129 and applied sciences. For instance, to study aggregation of naphthenic acids in aqueous
130 environments with different saline concentrations⁵⁰, in $\log D_{\text{pH}}$ calculations for lignin derivatives
131 and small datasets of drug-like compounds in different solvents by QM and statistical
132 thermodynamical methods⁵¹, partitioning of antioxidants⁵², aquatic hazard assessment of ionizable

133 organic chemicals⁵³, sorption mechanisms of antimicrobials in the soil⁵⁴, and physicochemical
134 properties of peptides and proteins.¹⁵⁻¹⁸

135 Previous studies have evaluated predictions of $\log D_{\text{pH}}$ using Equations 1 and 2 for a small
136 set of 35 ionizable molecules with computed $\log P_{\text{N}}$ and $\log P_{\text{I}}^{\text{app}}$ values calculated via an extension
137 of the Miertus-Scrocco-Tomassi solvation model.¹⁴ It was reported that Equation 1 tends to
138 overestimate the hydrophobicity of the studied molecules, given that the $P_{\text{I}}^{\text{app}}$ is not considered, on
139 the other hand, Equation 2 predicts a $\log D_{\text{pH}}$ value closer to the experimental values. This study
140 showed that Equation 2 provides a more exact lipophilicity profile at a wider pH range than
141 Equation 1. However, no systematized study has been performed to evaluate the importance of
142 considering the ionic partition on the $\log D_{\text{pH}}$ prediction for large sets of small drug-like molecules
143 at various pH values although when it has been reported that much of the poor performance of
144 some models on blind remains has been due to the simplification of ignoring the ionic species
145 partition.²⁷

146 In this study, we aim to provide guidance to the scientific community working in early-
147 stage drug design, food and environmental chemistry who deal with ionizable molecules, to
148 determine a priori which pH-dependent lipophilicity profile should be used depending on the
149 structure of a substance in their research. For this, we collected the experimental values of $\log P_{\text{N}}$,
150 $\text{p}K_{\text{a}}$, and $\log P_{\text{I}}^{\text{app}}$ of different compounds at various pH values, which are used to compute $\log D_{\text{pH}}$
151 with Equation 1 and Equation 2. We compared both calculations through statistical parameters
152 with the experimental $\log D_{\text{pH}}$ values. In addition, logistic regression (**LR**), random forest
153 classification (**RFC**), and support vector machine (**SVM**) models are developed to define from the
154 molecular structure which formalism is recommended for modeling a pH-dependent lipophilicity
155 profile.

156

157 **Methodology**

158 **Data collection and classification**

159 We critically compiled experimental values of $\log P_{\text{N}}$, $\text{p}K_{\text{a}}$, $\log P_{\text{I}}^{\text{app}}$, and $\log D_{\text{pH}}$ of 225
160 entries based on earlier literature reports (database available in reference 33).^{29,55,56} Refs. 29 and
161 55 were chosen based on the wide selection of experimental data for $\log P_{\text{N}}$, $\log D_{\text{pH}}$ and $\log P_{\text{I}}$

162 values and because they accommodate the desired chemical space of small molecules for our
163 modeling. SMILES codes were collected from publicly available data in PubChem⁵⁷ The pK_a
164 values were also obtained from PubChem but they were also corroborated by reviewing their
165 values in primary literature reports.^{38,57-80} The experimental technique of $\log P_N$, $\log D_{pH}$, and
166 $\log P_1^{app}$ measurements for each entry was thoroughly revised and added to the database.⁸¹⁻⁹¹ Ref 55
167 provided experimental $\log D_{pH}$ values of molecules at diverse pH ranges. The $\log P_1$ values were
168 obtained from the $\log D_{pH}$ at the most extreme measured pH, in which the molecule will be mostly
169 (above 95 %) in its ionized state. The $\log P_1^{app}$ values for molecules that were not measured under
170 ionizable pH conditions were obtained from external sources.^{38,74,92,93} The molecules were
171 classified as acids or bases based on their functional groups and pK_a values. Zwitterionic
172 compounds were found by evaluating the difference between acidic and basic pK_a in conjunction
173 with ChemAxon's calculator of protonated species distribution in function of pH.⁹⁴ Zwitterionic
174 and amphoteric species were also classified as acidic or basic based on their behavior of their
175 lipophilicity profiles, which were evaluated using the ChemAxon partitioning calculator.⁹⁵

176 The experimental data for each molecule were used to compute the $\log D_{pH}$ values using
177 Eq. 1 and Eq. 2 and are labeled as $\log D_{Eq,1}$ and $\log D_{Eq,2}$, respectively. The modeling performance
178 for each molecule was evaluated by calculating the absolute errors d_1 and d_2 (Eqs. 3 and 4):

$$179 \quad d_1 = \left| \log D_{Eq,1} - \log D_{exp} \right| \quad [3]$$

$$180 \quad d_2 = \left| \log D_{Eq,2} - \log D_{exp} \right| \quad [4]$$

181 where $\log D_{exp}$ represents the experimental $\log D_{pH}$ value.

182 The performances of the two formalisms were tested by performing a linear regression of
183 $\log D_{Eq,1}$ and $\log D_{Eq,2}$ on their experimental values. The root mean squared error (RMSE), mean
184 absolute error (MAE), mean squared error (MSE), and Pearson's correlation coefficient squared
185 (R^2) were calculated with the 'Metrics' package in R.⁹⁶ We also tested the performance of each
186 formalism on each individual molecule using descriptor d_3 (Eq. 5). When d_3 yields a value greater
187 than zero, Eq. 2 fits a more appropriate lipophilicity value and vice versa.

$$188 \quad d_3 = d_1 - d_2 \quad [5]$$

189

190 We create a binomial conditional based on the values of d_3 , where Eq. 2 should be used
191 when d_3 is greater than 0.2 (see Results and Discussion), otherwise, both equations are considered
192 to fit equally well, which can be interpreted as Eq. 1 providing better modeling due to its simplicity.

193

194 **Machine Learning models to classify the molecules according to the best fit to pH-dependent** 195 **lipophilicity profiles**

196 Topological and constitutional descriptors were calculated with the software ‘*rdck*’
197 package in R⁹⁷ while experimental descriptors ($\log P_N$, pK_a , and pH) were added from our dataset.
198 We also added the free energies of hydration and hydrogen bond strengths computed using the
199 open-source tool ‘*Jazzy*’⁹⁸ The H-bond donor and acceptor strengths were obtained by calculating
200 the partial charges of hydrogen atoms and atoms with lone electron pairs, respectively, along with
201 corrective terms. The free energy of hydration was calculated using the sum of the polar, apolar,
202 and interaction terms. The polar term was derived from the previously calculated H-bond donor
203 and acceptor strengths. The apolar terms consist of the sum of the weighted contributions of the
204 topological surface area, number of rings, and p-orbital counts in sp and sp² atoms. The interaction
205 term consists of a weighted sum of the amount of neighboring H-bond acceptor groups each atom
206 has in a molecule.⁹⁸

207 We eliminated intercorrelated properties so that no descriptor had a correlation value of r^2
208 > 0.6 (Figure S1 and S2). After this filtration step, two different feature selection methods were
209 tested to choose the best descriptors for our Machine Learning models. Firstly, we performed a
210 Welch’s *t*-test (**WTT**), which evaluates the statistical difference between the means of two
211 populations that have unequal variances and sample sizes.^{99,100} The algorithm calculates the mean
212 of both groups from the binomial conditional for each descriptor. These values are evaluated using
213 Equation 6.

$$214 \quad t = \frac{\Delta\mu}{\delta_{\Delta\bar{x}}} \quad [6]$$

215

216 where t stands for the statistic t in the Welch's t-test, $\Delta\mu$ represents the mean difference between
217 data samples from each population (Eq. 1 or Eq. 2 better fits), and the uncertainty value of both
218 groups, which was calculated using the standard deviation of both population samples (Eq.7):

219

$$220 \quad \delta_{\Delta\bar{x}} = \sqrt{\left(\frac{s_1}{\sqrt{N_1}}\right)^2 + \left(\frac{s_2}{\sqrt{N_2}}\right)^2} \quad [7]$$

221

222 The WTT was performed for each descriptor using R where the p -value was extracted.
223 Features that did not show statistical significance between the means ($p > 0.05$) were eliminated.
224 Secondly, a recursive feature elimination (**RFE**) was performed. This iterative feature selection
225 method builds a predictive model using the entire set of descriptors and calculates its importance
226 score (see Figure S3). The least important descriptors are removed, and the model was re-iterated
227 to achieve maximum performance.¹⁰¹ This RFE algorithm was programmed using the ‘*caret*’
228 package in R¹⁰² and tuned via a 5-time reiterated k -fold cross validation ($k = 10$). Table 1 shows
229 the selected descriptors with the WTT feature selection method for acids and bases, along with
230 their definitions and target molecules. Table S1 shows the descriptors selected using the RFE
231 method.

232

233

234

235

236

237

238

239

240

241

242 **Table 1.** List of the most influential structural descriptors^{98,103-106} used for the logistic regression
 243 models, their target molecules, and the divergence between the two populations from our dataset
 244 were determined using the WTT feature selection method by separating the populations with the
 245 conditional $d_3 > 0.2$.

Descriptor	Type	Definition	Target molecules
MDEC.11	Topological CDK descriptor	Molecular distance edge between all primary carbons.	Acids
MDEC.22		Molecular distance edge between all secondary carbons.	Acids
khs.sCH3		Number of -CH ₃ fragments in a molecule (Kier and Hall).	Acids
C2SP3		Singly bound carbon atom bound to two other carbons.	Acids
khs.dsCH		Number of =CH- fragments in a molecule (Kier and Hall).	Acids
khs.sNH2		Number of -NH ₂ fragments in a molecule (Kier and Hall).	Acids
khs.dssS		Number >S= fragments (sulfones) in a molecule (Kier and Hall).	Acids
HybRatio		Ratio of the number of sp^3 -C atoms compared to the sum of sp^3 and sp^2 C atoms.	Acids
C1SP3		Singly bound carbon atom bound to one other carbon.	Acids
nRings7		Number of 7-membered rings	Bases
khs.aaNH		Number of Ar-NH-Ar fragments in a molecule (Kier and Hall).	Bases
ATSc3		Autocorrelation topological distance weighed by charge calculated at every 3-atom distanced segment. Moreau-Broto autocorrelation descriptor 3 using polarizability	Bases
Alogp2	Constitutional CDK descriptor	$(\log P)^2$ value calculated with 3D structure directed QSAR method (Ghose & Grippen $\log K_{o/w}$).	Acids & Bases
delta	Experimental descriptor	δ (acids) = pH - pK_a δ (bases) = pK_a - pH	Acids & Bases
CH_strength	Jazzy calculation	C-H donor strength predicted with the Jazzy calculations.	Acids

246

247 **Logistic Regression Classification**

248 A logistic regression (LR) is a simple classification statistical model that provides a binary
249 response to the distribution of the input data among a specific descriptor. The simplest regressions
250 fit the distributions of data to a sigmoidal function, where the input values are given a probability
251 value, which is then classified into one of the two classes based on a cut-off value. We firstly
252 performed a feature selection process specific for logistic regressions by using the ‘*bestglm*’
253 package in R¹⁰⁷ which evaluates through n iterations, which combination of descriptors gives the
254 best fitted regression through the *leaps* algorithm.¹⁰⁸ This package evaluates the weight of each
255 descriptor by linearizing the sigmoidal function and giving a slope value and standard error for
256 each parameter like a multiple linear regression model (Equation 8).

$$257 \quad \ln\left(\frac{f(x)}{1-f(x)}\right) = \sum_{i=1}^n c_i x_i + b \quad [8]$$

258 The ‘*bestglm*’ package drops the parameters, where $c_i \rightarrow 0$. The algorithm iterates the
259 sigmoidal fit using Equation 8 n times until it finds the combination of descriptors in which the
260 parameters have the smallest standard error. This feature selection process was performed
261 separately for acids and bases because the descriptors have different behaviors for each type of
262 molecule.

263 Figure 2 shows a flowchart of the modelling process. The dataset was divided into acids
264 (113 entries) and bases (100 entries). Zwitterions (7 entries) were not considered for the Machine
265 Learning predictions because of their small sample size and because further lipophilicity modeling
266 can be performed for these molecules (see Results and Discussion section). Acids and bases were
267 randomly sampled into training and test sets at a ratio of 80:20. Multiple logistic regressions were
268 performed for the training sets based on previously collected descriptors. Predictive models were
269 programmed using the ‘*caret*’ package. Acids and bases were modeled separately and labeled as
270 **Models A** and **B**, respectively (see Figure 2). The test sets were evaluated using both models. The
271 performance of Models A and B was evaluated using confusion matrices (see Table S2), which
272 are widely used to evaluate classification models.¹⁰⁹ The confusion matrices tabulate the number
273 of true positives (**TP**), false positives (**FP**), true negatives (**TN**), and false negative (**FN**)
274 predictions, along with the sensitivity, specificity, and accuracy of the models. Sensitivity
275 determines the ability of the model to detect events of the positive class, that is, it indicates the
276 predictive performance of the molecules of the $\log D_{\text{Eq.2}}$ population (Equation 9). On the other

277 hand, specificity indicates the performance of the model in detecting the negative class, which in
278 this case are the molecules of the $\log D_{\text{Eq.1}}$ population (Equation 10). The accuracy indicates the
279 overall performance in detecting false positives and false negatives (Equation 11).

$$280 \quad \text{Sensitivity} = \frac{TP}{TP + FN} \quad [9]$$

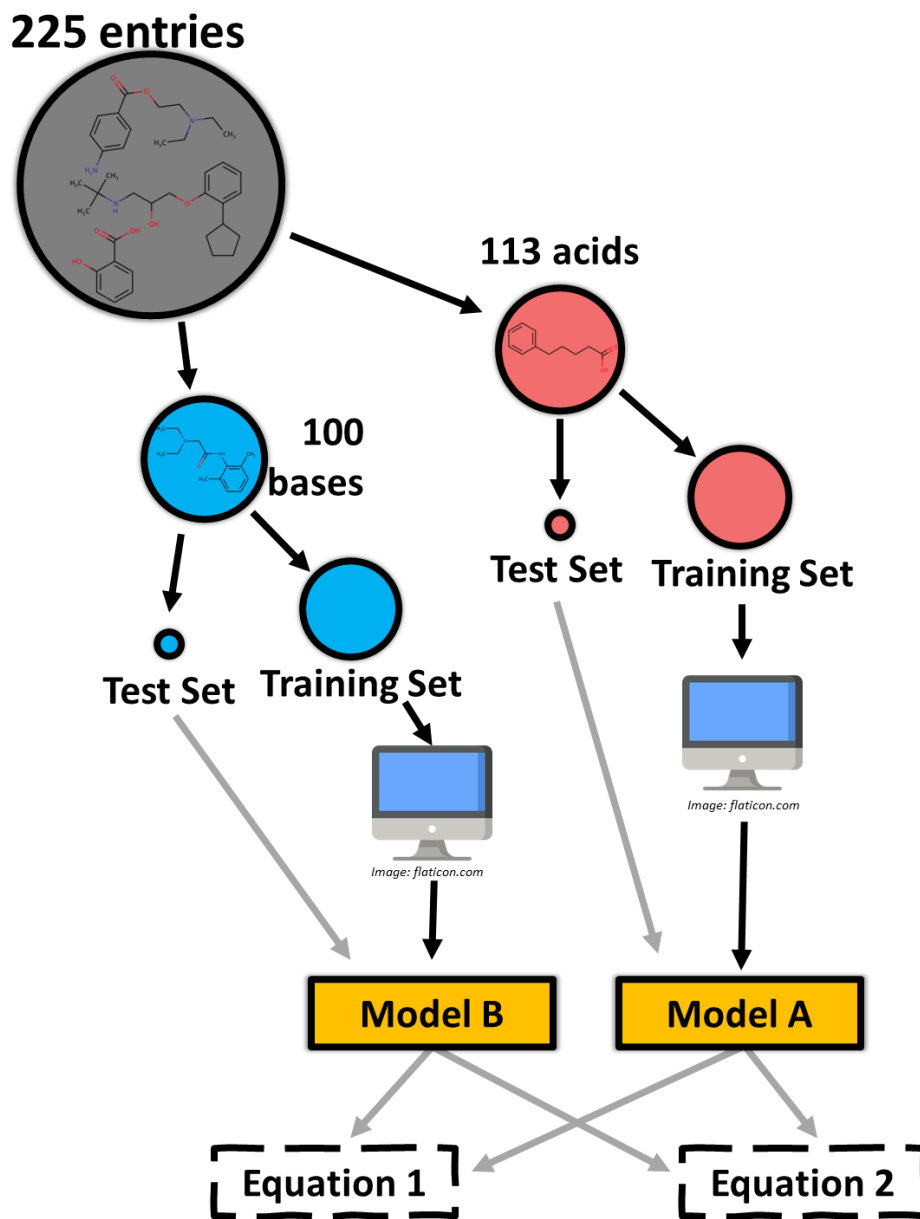
$$281 \quad \text{Specificity} = \frac{TN}{FP + TN} \quad [10]$$

$$282 \quad \text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad [11]$$

283

284 Models A and B were tested further using an external set. The experimental lipophilicity
285 measurements made by Disdier *et al.*⁴⁵ consisted of 69 data entries of small molecules with 38
286 acids, 16 bases, and 15 zwitterions, the latter being discarded for our analysis. To further check
287 the robustness of our models, a second external set of amino acid analogs were evaluated¹¹⁰,
288 consisting of 8 entries of histidine (basic amino acid) and 10 entries of tyrosine (acidic amino acid).
289 Then, we evaluated the performance of Model A and Model B for this data set with confusion
290 matrices (see Table S3-S4).

291



292

293 **Figure 2.** Graphical representation of the data classification and sampling of our dataset to create
 294 our predictive multiple logistic regression model using topological, constitutional, and
 295 experimental descriptors.

296

297

298 **Random Forest Classification**

299 Decision trees are a simple visual method for evaluating or classifying data, where each
300 node consists of a variable in the dataset. Each node leads to a leaf in which the desired output is
301 issued. A random forest is a combination of decision trees, which are randomly sampled, and the
302 nodes are randomly organized.¹¹¹ We split our dataset, training-, and test-sets as shown in Figure
303 2. In this case, Model A and Model B consist of random forest classification (**RFC**) performed
304 with the ‘*randomForest*’ package in R.¹¹² Both models were previously refined using the *tuneRF*
305 function within the package, which chooses the optimal *mtry* variable. This value indicates the
306 number of features selected at each split in each decision tree, where *mtry* = 2 gave the best
307 prediction for both models (number of trees = 500, see Supporting Information Figure S4). The
308 importance of each descriptor in both models was evaluated through the mean decrease in the Gini
309 impurity index using the *MeanDecreaseGini* function (see Figure S4).

310 The best lipophilicity profile fit for the acidic and basic tests and external sets was predicted with
311 Models A and B, respectively. The performance of each prediction was evaluated using confusion
312 matrices (see Tables S5-S7) and their respective sensitivity, specificity, and accuracy calculations
313 (Eqs. 9-11).

314

315 **Support Vector Machine Classification**

316 A Support Vector Machine (**SVM**) algorithm works by dividing training data into two
317 categories, either by linear or nonlinear classification; new data are then assigned to one of the two
318 classes. The model separates the data by finding a hyperplane that maximizes the gap between
319 categories. In the case of linear classification, the space is two-dimensional, making the
320 hyperplane a linear function.¹¹³ When the data are not linearly separable, the algorithm performs
321 the kernel trick, which consists of increasing the dimensions of the data space. This results in the
322 hyperplane being able to be another function in the original space, such as radial or polynomial,
323 allowing to classify the data in different ways.¹¹⁴

324 We split our datasets in the same manner as with the other classification models and set
325 Model A and Model B as support vector machines given by the ‘*e1071*’ package in R.¹¹⁵ We
326 decided to compare the performance of using a linear kernel (**SVML**) and a polynomial kernel
327 (**SVMP**); radial kernels were not evaluated because our binary data do not follow a circular

328 separation by the hyperplane, so it does not adequately fit a radial kernel SVM classification. The
329 hyperparameter selection for each model was performed with the *trainControl* and *train* functions
330 from the 'caret' package, which executes a *k*-fold cross-validation (*k* = 10 was used), where
331 different values of the parameters were tested and selected, which resulted in the highest accuracy.
332 The best hyperparameters were the function's default parameters: *C* = 1 for SVML and for SVMP,
333 *C* = 1, *degree* = 3, *gamma* = 1, and *coef0* = 0. We calculated the accuracy, sensitivity, and
334 specificity of each model using Eq. 9-11, using the results from their respective confusion matrices
335 (see Tables S8-S13). We then compared the confusion matrices of the LR, RFC, SVML, and
336 SVMP models to determine the one that yielded the best results.

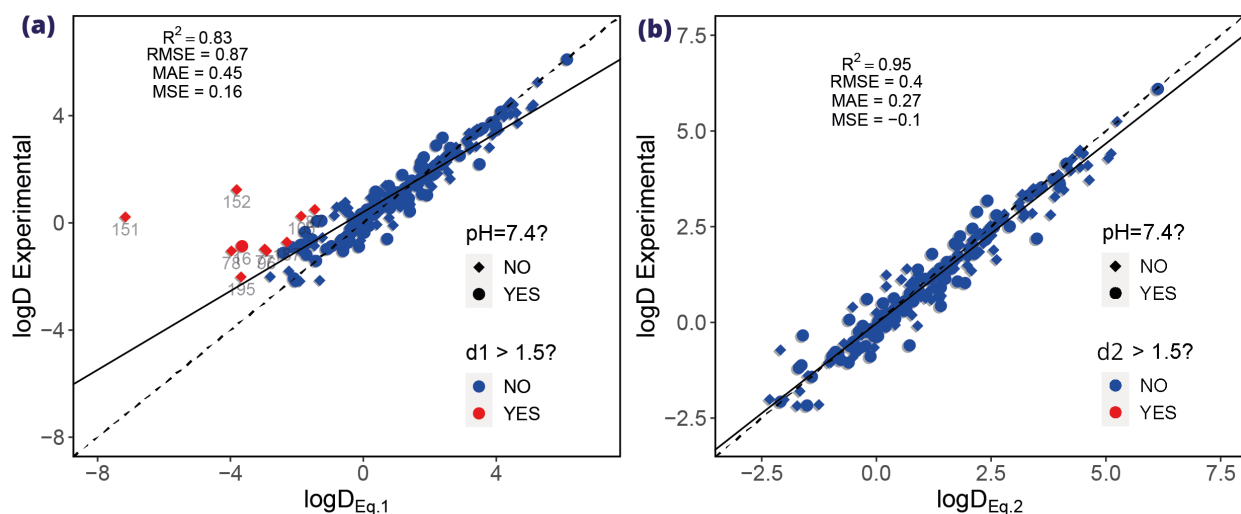
337

338 **Results and Discussion**

339 Our database consists of pK_a , $\log P_N$, $\log P_I^{app}$, and $\log D_{7.4}$ values reported by Avdeef²⁹. In
340 addition, we employed experimental entries of 86 molecules from the work of Tsantili-Kakoulidou
341 and collaborators containing $\log D_{pH}$ values at various pH for each molecule as an individual
342 entry.⁵⁵ Molecules with $\log D_{pH}$ values measured in the presence of background salt concentrations
343 above 0.15 mol/L were discarded because the study of the effect of external ions on lipophilicity
344 is beyond the scope of our study. Thus, we finally obtained 225 entries (118 individual molecules)
345 with 113 acids, 100 bases, and 12 zwitterions.

346 Calculation of $\log D_{pH}$ was accomplished using Eq. 1 and Eq. 2 for each molecule at their
347 respective pH. Figure 3 shows the overall performance of each model by comparing the computed
348 values with their respective experimental $\log D_{pH}$ values. As expected, most of the molecules
349 whose $\log D_{pH}$ values were measured under different pH conditions to 7.4, present the largest
350 deviation using the Eq. 1 (see Figure 3a, red marks) with highly underestimated predictions. As a
351 consequence, Eq. 1 poorly predicts the $\log D_{pH}$ values at extreme pH. On the other hand, the
352 predicted values using the Eq. 2 are significantly better (see Figure 3b), reducing the RMSE by
353 0.48 $\log D$ units, which represents an improvement of 55% in accuracy.

354



355
 356 **Figure 3.** Evaluation of the computed $\log D_{\text{pH}}$ of our database compared with the experimental
 357 values with (a) Eq. 1 and (b) Eq. 2. Rhomboids represent $\log D_{\text{pH}}$ when the pH is different of 7.4.
 358 Red dots and rhomboids highlight compounds with a deviation greater than 1.5 $\log D$ units.
 359 Statistical parameters were calculated using the ‘Metrics’ package in R (R^2 = squared Pearson’s
 360 correlation coefficient, RMSE = root mean squared error, MAE = mean absolute error and , MSE
 361 = mean squared error).

362
 363 Table 2 shows the reduction of RMSE in $\log D$ units of each molecule type by using Eq.2
 364 instead Eq.1. It is observed that any type of molecule shows a significant improvement in its
 365 performance when its distribution coefficient is modeled with $\log D_{\text{Eq.2}}$ (see Figure S5). Basic
 366 molecules showed the greatest improvement as the deviation shown by $\log D_{\text{Eq.1}}$ was greater than
 367 one unit of RMSE in $\log D$ units. Zwitterions also showed a significant improvement, even though
 368 these molecules can have multiple ionic partition coefficients (cationic partitions P^+ , and anionic
 369 partitions P^- , and zwitterionic partitions P_z), which are not considered in the model $\log D_{\text{Eq.2}}$. These
 370 partitions can be added by considering both acidic and basic pK_a into the thermodynamic
 371 equilibria.⁴⁵ Despite this, the implementation of just one of the two P_1^{app} did a significant
 372 improvement in the lipophilic modelling of zwitterions.

373
 374

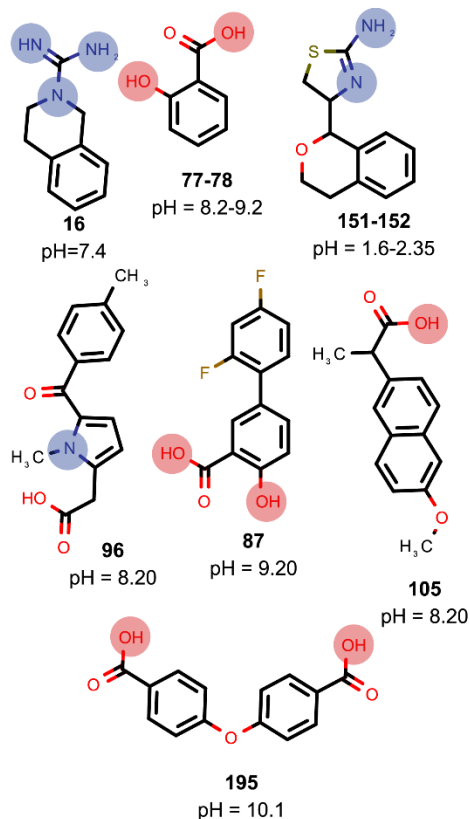
375 **Table 2.** Values of Δ RMSE for each type of molecule analyzed within our dataset by comparing
376 the modelled lipophilicities by $\log D_{\text{Eq.1}}$ and $\log D_{\text{Eq.2}}$ with their experimental values (Figure S1).

Type	Δ RMSE ^a
Acid	0.30
Base	0.67
Zwitterion	0.38
All	0.48

377 ^a Δ RMSE = RMSE ($\log D_{\text{Eq.1}}$) - RMSE ($\log D_{\text{Eq.2}}$)

378

379 The molecules with the highest deviations in the prediction of the experimental $\log D_{\text{pH}}$
380 using the $\log D_{\text{Eq.1}}$ are displayed in Figure 4. The chemical nature of the outliers is dominated by
381 the presence of ionic species because these compounds were experimentally measured to extreme
382 pH. These deviations respond to the theoretical framework of Eq. 1 and Eq. 2, thus, the inclusion
383 of the term P_1^{app} in Eq. 2 corrects the prediction. Figure 4 shows various polyacids or amphoteric
384 molecules with multiple ionizable sites included in our dataset. Bases **16**, **151-152** have multiple
385 protonation sites, while acids **77**, **78**, **87**, and **195** have two deprotonation sites, and amphoteric **96**
386 has a carboxylic acid and a tertiary aromatic nitrogen that can protonate at certain pH values. The
387 prediction of the $\log D_{\text{pH}}$ of these molecules can be improved by using more complex
388 thermodynamic models considering several equilibria.⁴⁵ However, it is shown here that the
389 consideration of one of the P_1^{app} with $\log D_{\text{Eq.2}}$ is enough to significantly increase the accuracy of
390 the lipophilicity modeling of these compounds to extreme pH where one charged species can
391 predominate over the others.



392

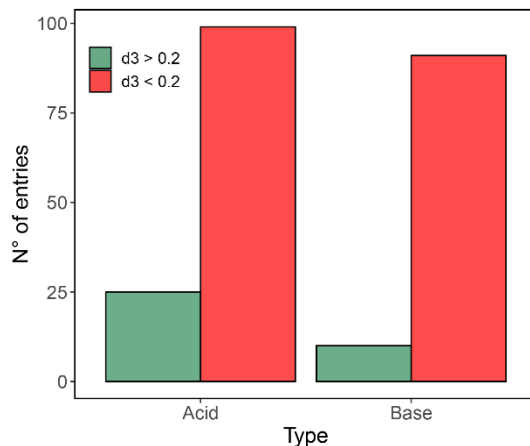
393 **Figure 4.** Representation of the molecules with the highest deviations in the prediction of the
 394 experimental $\log D_{\text{pH}}$ using the $\log D_{\text{Eq.1}}$. The protonation and deprotonation sites of each molecule
 395 were labeled in blue and red, respectively.

396

397 One of the aims of this study is to develop a classification algorithm that can differentiate
 398 whether the lipophilicity profile of a molecule would be better predicted with $\log D_{\text{Eq.1}}$ or $\log D_{\text{Eq.2}}$.
 399 However, we noticed that a significant number of entries yielded d_3 values close to 0 (see Figure
 400 S6a), which denotes that both formalisms compute a similar result compared to their experimental
 401 values. Therefore, let us note that we focus on the specific cases with a significant improvement
 402 when the P_1^{app} of molecules is considered. Indeed, we decided to delimit the conditional d_3
 403 indicating that if a molecule exceeds a certain value of d_3 , it is important to consider its apparent
 404 ionic partitioning for predicting its lipophilicity. We tested d_3 values between 0.1-1 and picked the
 405 optimal value based on two parameters. Firstly, considering that our set is small due to the fact
 406 that we used strictly experimental values in our database, we seek that the population of molecules
 407 that best fit with $\log D_{\text{Eq.2}}$ should be at least 10 %. Then, there should be a sufficient number of

408 descriptors that have statistically proven divergence by WTT ($p < 0.05$). Thus, machine learning
409 algorithms will have a larger number of parameters to create predictive models with higher
410 accuracy. In consequence, the delimiter '0.2' showed an adequate balance between these two
411 parameters, and it was selected as our cut-off value (see Figure S6b). Thus, molecules with values
412 of $d_3 > 0.2$ showed an improvement in lipophilicity modeling using Eq. 2. On the other hand,
413 entries that had negative d_3 values or that fell into the range $0.2 < d_3 < 0$ were classified as
414 molecules where the difference between both models was negligible, and thus were classified as
415 better fitted using the $\log D_{\text{Eq.1}}$ due to its easy implementation (it does not depend on P_1^{app} , resulting
416 in less computational effort and fewer experimental parameters). Higher thresholds significantly
417 decreased the population in $\log D_{\text{Eq.2}}$, while lower values reduced the structural divergence between
418 molecules in $\log D_{\text{Eq.1}}$ and $\log D_{\text{Eq.2}}$, making it more difficult to find structural descriptors that can
419 differentiate between both populations. The value '0.5' was also tested since a local maximum of
420 descriptors with p-values < 0.05 was observed at this point (see Figure S6b). Furthermore, this
421 value is of experimental interest, because $\log P_{\text{N}}$ measurements of substances with different
422 techniques tend to vary by amounts less than $0.5 \log P$ units (using the Shake-Flask method as a
423 reference), being this value considered as a parameter to indicate that the experimental techniques
424 are not equivalent.³⁰ However, this extreme value and the descriptors selected (see Table S14)
425 showed poor performance in the ML models tested, especially with the external set 1 (see Figure
426 S7). This phenomenon can be explained since this d_3 delimiter has a very small $\log D_{\text{Eq.2}}$
427 population, thus the datasets are extremely unbalanced, and the robustness of the models is
428 reduced, on the other hand, the accuracy of experimental methods, even using different techniques,
429 rounds at values less than $0.2 \log P$ units.³⁰ Therefore, we continued to train the ML models using
430 the $d_3 > 0.2$ cut-off value to determine tendencies among the selected descriptors via the feature
431 selection methods and to evaluate the performance of the ML algorithms.

432 Figure 5 shows the distribution of the molecules in our database, classified using the criteria
433 $d_3 > 0.2$ as binary descriptor. Most entries can be computed using $\log D_{\text{Eq.1}}$ with satisfactory results.
434 However, we observed that 25 acids and 10 bases showed a clear improvement within our d_3
435 threshold by modeling lipophilicity with $\log D_{\text{Eq.2}}$.



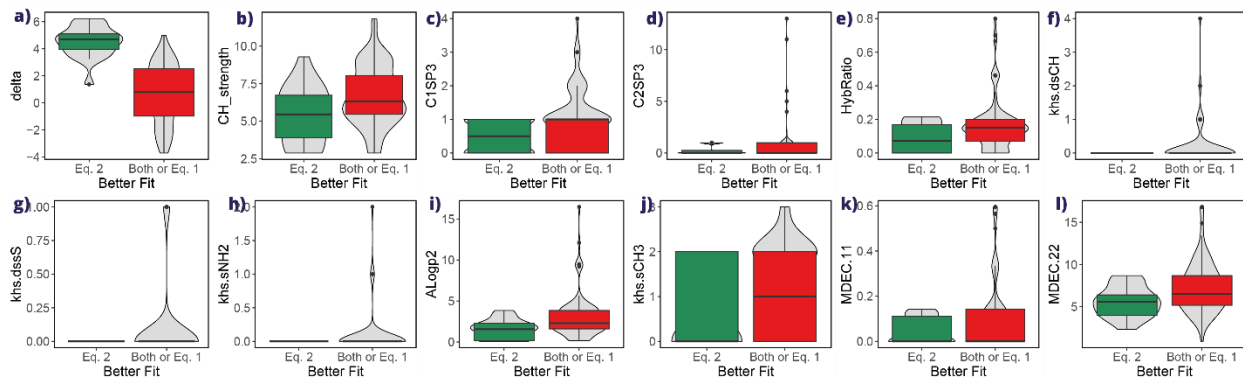
436
 437 **Figure 5.** Distribution of acid and basic entries from our dataset as a function of their d_3 values.

438
 439 We obtained several structural and physicochemical descriptors of the molecules to find a
 440 considerable divergence between the populations. First, our database was split into acids and bases
 441 and then in training and test set. The ‘*rcdk*’ package in R was used to look through the descriptors,
 442 along with the *Jazzy* calculations of energies of hydration and hydrogen-bond strengths and the
 443 experimental descriptors. The feature selection methods selected show a wide range of diverse
 444 descriptors (see Table 1 and Table S1). Then, we performed a Welch’s *t*-test on our descriptors
 445 (WTT) where is analyzed the divergence between populations relative to the variances of the two
 446 groups.⁹⁹ This test was selected over a Student’s *t*-test because of the divergence of sample sizes
 447 (Figure 5) and variances between groups (Figure 6-7).¹⁰⁰ The WTT descriptors gave acceptable
 448 accuracies (see Figure S8).

449 An iterative feature selection method was also tested using an RFE model. The algorithm
 450 achieved better performance when the 14 most important variables for acids and the nine most
 451 important variables for bases were maintained. The importance of each descriptor posed by RFE
 452 is shown in Figure S3. Good results were obtained when these descriptors were implemented in
 453 the training of the machine learning models. However, the accuracy decreased significantly when
 454 the test and external set 1 was evaluated (see Figure S8c-d), indicating that these descriptors did
 455 not generate a sufficiently robust model, or a large number of chosen descriptors (see Table S1)
 456 may overfit the data. Therefore, we selected the WTT descriptors to analyze the tendencies of the

457 molecules in each population and to evaluate the overall performance of the machine learning
458 algorithms that we developed.

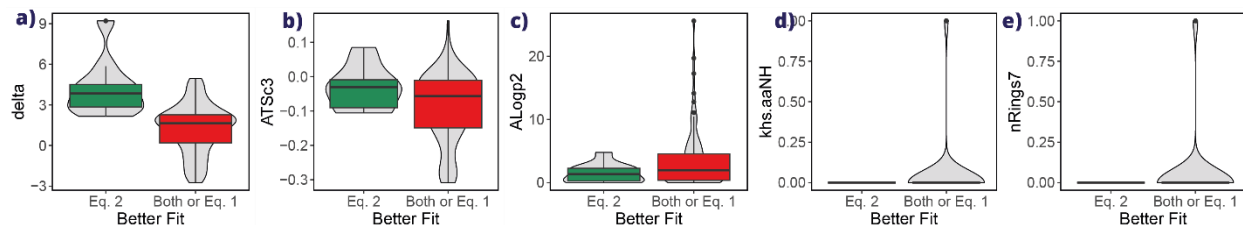
459



460

461 **Figure 6.** Violin plots of the distribution of the acidic molecules in our dataset along the selected
462 descriptors for the acids ((a) *delta*, (b) *CH_strength*, (c) *C1SP3*, (d) *C2SP3*, (e) *HybRatio*, (f)
463 *khs.dsCH*, (g) *khs.dssS*, (h) *khs.sNH2*, (i) *Alogp2*, (j) *khs.sCH3*, (k) *MDEC.11*, and (l) *MDEC.22*).
464 Distributions are separated between acids and bases and classified by the binary operator $d_3 > 0.2$
465 (green) and $d_3 < 0.2$ (red).

466



467

468

469 **Figure 7.** Violin plots of the distribution of the acidic molecules in our dataset along the selected
470 descriptors for the bases (a) *delta*, (b) *ATSc3*, (c) *Alogp2*, (d) *khs.aanNH*, and (e) *nRings7*).
471 Distributions are separated between acids and bases and classified by the binary operator $d_3 > 0.2$
472 (green) and $d_3 < 0.2$ (red).

473

474

475 Figure 6 and Figure 7 show the selected descriptors for acids and bases used to train our
476 classification ML models, respectively. These descriptors showed a statistically significant
477 divergence between the means of both populations among 180 descriptors tested for acids and
478 bases.

479 Both, acidic and basic compounds show significant differences in their means ($p < 0.005$
480 in WTT test) for the descriptors δ and $Alogp2$ (see Table 1). The descriptor δ was calculated
481 at the respective pH of each entry for the acids and bases. As expected, this descriptor correlates
482 with the prominence of ionic species in both phases. Therefore, the apparent ionic partition became
483 more significant for entries with higher δ values (see Figures 6a and Figure 7a). This result is
484 very promising, because despite being an experimental descriptor, there are computational
485 methods to determine the pK_a that include first principles models¹¹⁶⁻¹¹⁹ as well as machine learning
486 tools^{120,121}, so the descriptor δ can be automated and easily used to classify molecules
487 according to the lipophilicity formalisms analyzed here. In fact, the root-mean-square error
488 (RMSE) between predicted pK_a values using the software ChemAxon and experimental data in
489 our database is just 0.58 log units and the squared coefficient of determination (R^2) of 0.95 (see
490 Fig. S9)

491 The $Alogp2$ descriptor consists of a 3D-QSAR model by Ghose & Crippen (1986) that
492 predicts a square value of $\log P_N$ value by analyzing the presence of 110 structural fragments
493 within the molecules.¹⁰⁴ Figure 6i and 7c show that molecules with hydrophobicity close to $\log P$
494 = 0 (with lower $Alogp2$ values) tend to fit best with $\log D_{Eq.2}$. Water and *n*-octanol are not miscible,
495 yet a small amount of water can dissolve in octanol at room temperature (~ 2.9 mol/kg).¹⁰⁵ These
496 hydrophilic molecules might be dragged by the dissolved water to the octanol phase along with
497 ionic species; thus, the apparent ionic partition would have a higher importance in these molecules.

498

499 This affinity for water, at least for acidic compounds, was further demonstrated by the
500 $CH_strength$ descriptor (Figure 6b). This descriptor, calculated by *Jazzy*, predicts the hydrogen-
501 bond donor strength in carbon atoms.⁹⁸ The smaller $CH_strength$ values indicate that for entries
502 with $d_3 < 0.2$, H-bond donors are not primarily found on carbons. Instead, they are found on other
503 more electronegative heteroatoms. Thus, by weakening the X-H covalent bonds through H-bonds,
504 the possibility of ionization of these species in both water and *n*-octanol increases. Figure 6 present

505 other important descriptor for acidic compounds such as *MDEC.22* and *HypRatio*. The *MDEC.22*
506 descriptor consists of a relationship between the number of secondary carbons in the molecule
507 (i.e., vertices in a graph with only two paths) and the squared average atomic distance between
508 those atoms.¹⁰¹, whereas *HypRatio* is the number of sp³-C atoms compared to the sum of sp³ and
509 sp² C atoms. Eq. 2 works better for acidic substances with low values of these descriptors, which
510 considering together the values of *Alogp2*, allows us to intuit that small and rigid ionizable
511 molecules with instaurations or aromatic systems need considering the P_1^{app} to obtain an accurate
512 prediction of $\log D_{\text{pH}}$.

513 Similarly, for basic compounds, higher values of *ATSc3* descriptor are associated with
514 taking into account the P_1^{app} for modeling pH-dependent lipophilic profiles. This descriptor is
515 related with high molecular polarizability which agrees with the pattern of small molecules with
516 the presence of polar atoms such as nitrogen.

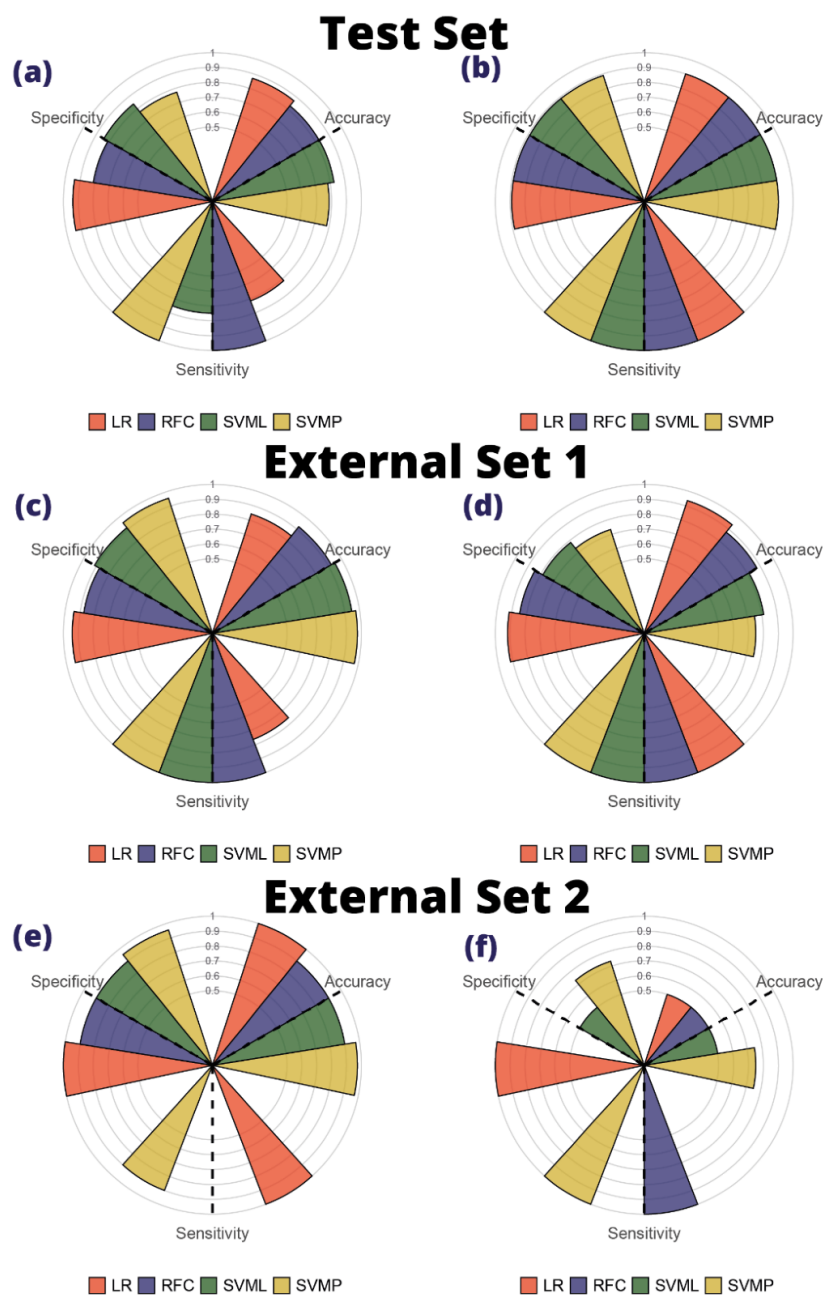
517 Therefore, the apparent ionic partition effect should be considered for these small, rigid,
518 and unsaturated molecules which present a significant proportion of ionic species in the aqueous
519 phase species. It has been previously shown that the P_1^{app} of molecules may mechanistically occur
520 via a simple ion-transfer reaction.¹²² Thus, it is more plausible that small and compact molecules
521 have a more prominent P_1^{app} because of the lower energetic cost of transferring to the cavity of the
522 ion they replace.

523

524

525 **Machine Learning Classification Models**

526 Models A and B (see Figure 2) were trained using the LR, RFC, and SVM algorithms. A training
527 set for acidic and basic molecules was used for each model and evaluated using the test set
528 consisting of 20% of our population (see Figure S10). In addition, two external sets were validated
529 with the experimental data of Disdier et al. (external set 1)⁴⁵ and Fauchère and Pliška (external set
530 2)¹¹⁰. Predictions were made as to which formalism best modeled the lipophilicity of the inputs,
531 and the results were collected in confusion matrices. The performance of each marker was
532 evaluated by calculating its accuracy, specificity, and sensitivity. Figure 8 shows the results of the
533 calculations for the four algorithms for the test and external sets of acidic and basic molecules.



534

535

536 **Figure 8.** Accuracy, sensitivity, and specificity of every ML model evaluated in this study for
 537 acidic (a,c,e) and basic (b,d,f) entries within the test and external sets by defining our populations
 538 with the conditional $d_3 > 0.2$. Descriptors were selected with the WTT method. Accuracies,
 539 sensitivities, and specificities were calculated with Eqs. 9-11 based on the results of each confusion
 540 matrix (Tables S1-S8)

541 It is observed that most of the calculated accuracies have high values (between 0.8 and
542 0.95), denoting that these classification models manage to distinguish relatively well which
543 molecules best fit with $\log D_{\text{Eq.1}}$ and $\log D_{\text{Eq.2}}$. However, it was observed that in the test set of acidic
544 molecules, the sensitivity decreased, indicating that the models had difficulties in detecting
545 molecules that fit $\log D_{\text{Eq.2}}$ (Figure 7a). The external set related with capping amino acids reported
546 by Fauchère and Pliška¹¹⁰ obtained divergent results. On the one hand, the pH-dependent values
547 of N-Acetyl-L-tyrosine amide were predicted with excellent metrics, especially using the LR and
548 SVM models, because our training set had a representative amount of phenolic groups. On the
549 other hand, in the case of N-Acetyl-L-histidine amide, the results were very poor, this is due, at
550 least in part, to the fact that our set has few bases in relation to the acids that best-fit to Eq. 2, and
551 mainly because there was no imidazole fragment present in our set of bases, thus limiting the
552 performance of our models.

553
554

555 **Conclusions**

556 Lipophilicity is undoubtedly the most used and important descriptor in the early stages of
557 drug discovery and development. Additionally, it is a crucial descriptor in substance risk
558 assessment and also in areas including adsorption in materials, catalysts, food chemistry, and
559 computational biology. There are multiple tools to determine this descriptor, mainly for neutral
560 molecules ($\log P_{\text{N}}$), and for substances with ionizable groups, two formalisms are commonly used
561 to determine the distribution coefficient ($\log D_{\text{pH}}$), being the simplest pH correction model the most
562 widely used. However, previous studies carried out on specific and small molecule sets
563 recommend considering the effect of the apparent ionic compounds (P_1^{app}), since it has seen a
564 negative impact on the accuracy of computing lipophilic profiles when charged species or related
565 species are ignored. Our study, which was based on a larger amount of data and strictly on
566 experimental values, validates the observations presented in limited previous studies. Thus, we
567 develop machine learning algorithms using logistic regressions, random forest classifications, and
568 support vector machine models to determine from molecular structures in which cases the P_1^{app}
569 should be considered. The results indicate that small, rigid, and unsaturated molecules with $\log P_{\text{N}}$

570 close to zero which present a significant proportion of ionic species in the aqueous phase, are better
571 modeled using the formalism which takes into account the apparent ionic compounds (P_1^{app}).

572 Although we are aware of the molecular complexity of the species that can be included for
573 the computational determination of the apparent ionic partition (P_1^{app}), parameterization or training
574 of models using experimental values of P_1^{app} can help to alleviate the restricted application of
575 formalisms that include this effect. Finally, our findings can serve as guidance to the scientific
576 community working in early-stage drug design, food, and environmental chemistry who deal with
577 ionizable molecules, to determine a priori which lipophilicity profile should be used depending on
578 the structure of a substance in research efforts. Future studies will address the influence played by
579 the apparent ionic partition (P_1^{app}) on the pH-dependent lipophilic profiles in more complex
580 systems such as zwitterionic and peptides.

581 **References**

582 (1) Waring, M. J. Lipophilicity in Drug Discovery. *Expert Opin Drug Discov* **2010**, *5* (3), 235–
583 248. <https://doi.org/10.1517/17460441003605098>.

584 (2) Lewis, D. F. V.; Jacobs, M. N.; Dickins, M. Compound Lipophilicity for Substrate Binding
585 to Human P450s in Drug Metabolism. *Drug Discov Today* **2004**, *9* (12), 530–537.
586 [https://doi.org/https://doi.org/10.1016/S1359-6446\(04\)03115-0](https://doi.org/https://doi.org/10.1016/S1359-6446(04)03115-0).

587 (3) Chatzopoulou, M.; Emer, E.; Lecci, C.; Rowley, J. A.; Casagrande, A. S.; Moir, L.; Squire,
588 S. E.; Davies, S. G.; Harriman, S.; Wynne, G. M.; Wilson, F. X.; Davies, K. E.; Russell, A. J.
589 Decreasing HepG2 Cytotoxicity by Lowering the Lipophilicity of Benzo[d]Oxazolephosphinate
590 Ester Utrophin Modulators. *ACS Med Chem Lett* **2020**, *11* (12), 2421–2427.
591 <https://doi.org/10.1021/ACSMEDCHEMLETT.0C00405/ASSET/IMAGES/LARGE/ML0C0040>
592 [5_0003.JPEG](https://doi.org/10.1021/ACSMEDCHEMLETT.0C00405/ASSET/IMAGES/LARGE/ML0C00405_0003.JPEG).

593 (4) Miller, M. M.; Wasik, S. P.; Huang, G. L.; Shlu, W. Y.; Mackay, D. Relationships between
594 Octanol-Water Partition Coefficient and Aqueous Solubility. *Environ Sci Technol* **1985**, *19* (6),
595 522–529. https://doi.org/10.1021/ES00136A007/ASSET/ES00136A007.FP.PNG_V03.

596 (5) Soliman, K.; Grimm, F.; Wurm, C. A.; Egner, A. Predicting the Membrane Permeability
597 of Organic Fluorescent Probes by the Deep Neural Network Based Lipophilicity Descriptor
598 DeepFl-LogP. *Scientific Reports 2021 11:1* **2021**, *11* (1), 1–9. [https://doi.org/10.1038/S41598-](https://doi.org/10.1038/S41598-021-86460-3)
599 [021-86460-3](https://doi.org/10.1038/S41598-021-86460-3).

600 (6) Esser, H. O. A Review of the Correlation between Physicochemical Properties and
601 Bioaccumulation. *Pestic Sci* **1986**, *17* (3), 265–276. <https://doi.org/10.1002/PS.2780170310>.

602 (7) Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and
603 Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and

- 604 Development Settings. *Adv Drug Deliv Rev* **2001**, *46* (1–3), 3–26. [https://doi.org/10.1016/S0169-409X\(00\)00129-0](https://doi.org/10.1016/S0169-409X(00)00129-0).
- 606 (8) Veber, D. F.; Johnson, S. R.; Cheng, H. Y.; Smith, B. R.; Ward, K. W.; Kopple, K. D.
607 Molecular Properties That Influence the Oral Bioavailability of Drug Candidates. *J Med Chem*
608 **2002**, *45* (12), 2615–2623.
609 https://doi.org/10.1021/JM020017N/SUPPL_FILE/JM020017N_S.PDF.
- 610 (9) Leo, A.; Hansch, C.; Church, C. Comparison of Parameters Currently Used in the Study of
611 Structure-Activity Relationships. *J Med Chem* **1969**, *12* (5), 766–771.
612 https://doi.org/10.1021/JM00305A010/ASSET/JM00305A010.FP.PNG_V03.
- 613 (10) Leo, A.; Hansch, C.; Elkins, D. Partition Coefficients and Their Uses. *Chem Rev* **1971**, *71*
614 (6), 525–616. https://doi.org/10.1021/CR60274A001/ASSET/CR60274A001.FP.PNG_V03.
- 615 (11) Kerns, E. H. High Throughput Physicochemical Profiling for Drug Discovery. *J Pharm Sci*
616 **2001**, *90* (11), 1838–1858. <https://doi.org/10.1002/JPS.1134>.
- 617 (12) Liu, X.; Tu, M.; Kelly, R. S.; Chen, C.; Smith, B. J. Development of a Computational
618 Approach to Predict Blood-Brain Barrier Permeability. *Drug Metab Dispos* **2004**, *32* (1), 132–
619 139. <https://doi.org/10.1124/DMD.32.1.132>.
- 620 (13) Colmenarejo, G. In Silico Prediction of Drug-Binding Strengths to Human Serum Albumin. *Med*
621 *Res Rev* **2003**, *23* (3), 275–301. <https://doi.org/10.1002/MED.10039>.(14)
- 622 (14) Zamora, W. J.; Curutchet, C.; Campanera, J. M.; Luque, F. J. Prediction of pH-Dependent
623 Hydrophobic Profiles of Small Molecules from Miertus-Scrocco-Tomasi Continuum Solvation
624 Calculations. *Journal of Physical Chemistry B* **2017**, *121* (42), 9868–9880.
- 625 (15) Iglesias, V., Pintado-Grima, C., Santos, J., Forn, M., Ventura, S. Prediction of the Effect
626 of pH on the Aggregation and Conditional Folding of Intrinsically Disordered Proteins with
627 SolupHred and DispHred. In: Carugo, O., Eisenhaber, F. (eds) *Data Mining Techniques for the*
628 *Life Sciences. Methods in Molecular Biology*. **2022**. vol 2449. Humana, New York, NY.
629 https://doi.org/10.1007/978-1-0716-2095-3_8.
- 630 (16) Oeller M, Kang R, Bell R, Ausserwöger H, Sormanni P, Vendruscolo M. Sequence-based
631 prediction of pH-dependent protein solubility using CamSol. *Brief Bioinform*. **2023**,
632 *19*;24(2):bbad004. doi: 10.1093/bib/bbad004.
- 633 (17) Porto WF, Ferreira KCV, Ribeiro SM, Franco OL. Sense the moment: A highly sensitive
634 antimicrobial activity predictor based on hydrophobic moment. *Biochim Biophys Acta Gen Subj*.
635 *2022 Mar*;1866(3):130070. doi: 10.1016/j.bbagen.2021.130070.
- 636 (18) Zamora, W. J.; Pinheiro, S.; Separovic, F.; Luque, F. J. Insights into the Effect of the
637 Membrane Environment on the Three-dimensional Structure-function Relationship of
638 Antimicrobial Peptides. *Biophysical Journal*, 2020, *118*, 3, 236a
- 639 (19) Simm, S.; Einloft, J.; Mirus, O.; Schleiff, E. 50 Years of Amino Acid Hydrophobicity
640 Scales: Revisiting the Capacity for Peptide Classification. *Biol Res* **2016**, *49* (1), 1–19.
641 <https://doi.org/10.1186/S40659-016-0092-5>.
- 642 (20) Zamora, W. J.; Campanera, J. M.; Luque, F. J. Development of a Structure-Based, PH-
643 Dependent Lipophilicity Scale of Amino Acids from Continuum Solvation Calculations. *Journal*
644 *of Physical Chemistry Letters* **2019**, *10* (4), 883–889.

- 645 (21) Ingram, T., Richter, U., Mehling, T., & Smirnova, I. Modelling of pH-dependent n-
646 octanol/water partition coefficients of ionizable pharmaceuticals. *Fluid Phase Equilibria*. 2011,
647 305, 197-203.
- 648 (22) Chen, Chao and Shiang-Tai Lin. "Prediction of pH Effect on the Octanol–Water Partition
649 Coefficient of Ionizable Pharmaceuticals." *Industrial & Engineering Chemistry Research*.
650 2016, 55, 9284-9294.
- 651 (23) Westall, J. C.; Leuenberger, C.; Schwarzenbach, R. P. Influence of pH and Ionic Strength
652 on the Aqueous-Nonaqueous Distribution of Chlorinated Phenols. *Environ Sci Technol* **1985**, 19
653 (2), 193–198.
- 654 (24) Xing, L.; Glen, R. C. Novel Methods for the Prediction of LogP, pKa, and LogD. *J Chem*
655 *Inf Comput Sci* **2002**, 42 (4), 796–805.
- 656 (25) Tetko, I. V.; Poda, G. I. Application of ALOGPS 2.1 to Predict Log D Distribution
657 Coefficient for Pfizer Proprietary Compounds. *J Med Chem* **2004**, 47 (23), 5601–5604.
- 658 (26) Livingstone, D. Theoretical Property Predictions. *Curr Top Med Chem* **2005**, 3 (10), 1171–
659 1192. <https://doi.org/10.2174/1568026033452078>
- 660 (27) C. C. Bannan, G. Calabro, D. Y. Kyu and D. L. Mobley, *J. Chem. Theory Comput.*, 2016,
661 **12**, 4015–4024.
- 662 (28) Bergazin, T. D.; Tielker, N.; Zhang, Y.; Mao, J.; Gunner, M. R.; Francisco, K.; Ballatore,
663 C.; Kast, S. M.; Mobley, D. L. Evaluation of Log P, PK a, and Log D Predictions from the
664 SAMPL7 Blind Challenge. *J Comput Aided Mol Des* **2021**, 35 (7), 771–802.
665 <https://doi.org/10.1007/s10822-021-00397-3>.
- 666 (29) Avdeef, A. *Absorption and Drug Development.*; John Wiley & Sons, **2012**.
- 667 (30) Port, A.; Bordas, M.; Enrech, R.; Pascual, R.; Rosés, M.; Ràfols, C.; Subirats, X.; Bosch,
668 E. Critical Comparison of Shake-Flask, Potentiometric and Chromatographic Methods for
669 Lipophilicity Evaluation (Log P O/W) of Neutral, Acidic, Basic, Amphoteric, and Zwitterionic
670 Drugs. *European Journal of Pharmaceutical Sciences* **2018**, 122, 331–340.
671 DOI:10.1016/j.ejps.2018.07.010.
- 672 (31) Austin, Rupert P. et al. The effect of ionic strength on liposome-buffer and 1-octanol-buffer
673 distribution coefficients. *Journal of Pharmaceutical Sciences*. **1998**, 87, 5, 599-607 .
- 674 (32) Jain P, Kumar A. Concentration-dependent apparent partition coefficients of ionic liquids
675 possessing ethyl- and bi-sulphate anions. *Phys Chem Chem Phys*. **2016**, 14;18(2):1105-13. doi:
676 10.1039/c5cp06611e.
- 677 (33) Jafvert, C. T., Westall, J. C., Grieder, E. & Schwarzenbach, R. P. Distribution of
678 Hydrophobic Ionogenic Organic Compounds Between Octanol and Water: Organic Acids.
679 *Environ. Sci. Technol.* **1990**, **24**, 1795–1803.
- 680 (34) Austin, R. P., Davis, a. M. & Manners, C. N. Partitioning of Ionizing Molecules between
681 Aqueous Buffers and Phospholipid Vesicles. *J. Pharm. Sci.* **1995**, **84**, 1180–1183.
- 682 (35) Takács-Novák, K. & Szász, G. Ion-Pair Partition of Quaternary Ammonium Drugs: The
683 Influence of Counter Ions of Different Lipophilicity, Size, and Flexibility. *Pharm. Res.* **1999**, **16**,
684 1633–1638.

- 685 (36) Fini, A. et al. Formation of Ion-Pairs in Aqueous Solutions of Diclofenac Salts. *Int. J.*
686 *Pharm.* **1999**, **187**, 163–173.
- 687 (37) Sarveiya, V., Templeton, J. F. & Benson, H. a E. Ion-Pairs of Ibuprofen: Increased
688 Membrane Diffusion. *J. Pharm. Pharmacol.* **2004**, **56**, 717–724.
- 689 (38) Scherrer, R. A. & Donovan, S. F. Automated Potentiometric Titrations in KCl/Water-
690 Saturated Octanol: Method for Quantifying Factors Influencing Ion-Pair Partitioning. *Anal. Chem.*
691 **2009**, **81**, 2768–2778
- 692 (39) Wenlock, M. C., Potter, T., Barton, P. & Austin, R. P. A Method for Measuring the
693 Lipophilicity of Compounds in Mixtures of 10. *J. Biomol. Screen.* **2011**, **16**, 348–55.
- 694 (40) Fini, A., Bassini, G., Monastero, A. & Cavallari, C. Diclofenac Salts, VIII. Effect of the
695 Counterions on the Permeation through Porcine Membrane from Aqueous Saturated Solutions.
696 *Pharmaceutics.* **2012**, **4**, 413–429
- 697 (41) Paternostre M, Meyer O, Grabielle-Madelmont C, Lesieur S, Ghanam M, Ollivon M.
698 Partition coefficient of a surfactant between aggregates and solution: application to the micelle-
699 vesicle transition of egg phosphatidylcholine and octyl beta-D-glucopyranoside. *Biophys J.* 1995
700 Dec;69(6):2476-88. doi: 10.1016/S0006-3495(95)80118-9.
- 701 (42) Pieńko T, Grudzień M, Taciak PP, Mazurek AP. Cytisine basicity, solvation, logP, and
702 logD theoretical determination as tool for bioavailability prediction. *J Mol Graph Model.* **2016**
703 ;63:15-21. doi: 10.1016/j.jmgm.2015.11.003.
- 704 (43) Hung, L. Q. Electrochemical Properties of the Interface between Two Immiscible
705 Electrolyte Solutions. *J. Electroanal. Chem.* **1980**, **115**, 159–174.
- 706 (44) Kakiuchi, T. Limiting Behavior in Equilibrium Partitioning of Ionic Components in
707 Liquid–Liquid Two-Phase Systems. *Anal. Chem.* **1996**, **68**, 3658.
- 708 (45) Disdier, Z.; Savoye, S.; Dagnelie, R. V. H. Effect of Solutes Structure and PH on the N-
709 Octanol/Water Partition Coefficient of Ionizable Organic Compounds. *Chemosphere* **2022**, **304**.
710 <https://doi.org/10.1016/j.chemosphere.2022.135155>.
- 711 (46) Berthod, A.; Carda-Broch, S.; Garcia-Alvarez-Coque, M. C. Hydrophobicity of Ionizable
712 Compounds. A Theoretical Study and Measurements of Diuretic Octanol-Water Partition
713 Coefficients by Countercurrent Chromatography. *Anal Chem* **1999**, **71** (4), 879–888.
- 714 (47) Âde, F.; Reymond, R.; Carrupt, P.-A.; Testa, B.; Girault, H. H. Charge and Delocalisation
715 Effects on the Lipophilicity of Protonable Drugs. *Chemistry – A European Journal* **1999**, **5** (1),
716 39–48.
- 717 (48) Gobry, V.; Ulmeanu, S.; Reymond, F.; Bouchard, G.; Carrupt, P. A.; Testa, B.; Girault, H.
718 H. Generalization of Ionic Partition Diagrams to Lipophilic Compounds and to Biphasic Systems
719 with Variable Phase Volume Ratios. *J Am Chem Soc* **2001**, **123** (43), 10684–10690.
- 720 (49) Reymond, F.; Steyaert, G.; Carrupt, P. A.; Testa, B.; Girault, H. Ionic Partition Diagrams:
721 A Potential-PH Representation. *J Am Chem Soc* **1996**, **118** (47), 11951–11957.
- 722 (50) Cunha, R. D.; Ferreira, L. J.; Orestes, E.; Coutinho-Neto, M. D.; de Almeida, J. M.;
723 Carvalho, R. M.; Maciel, C. D.; Curutchet, C.; Homem-de-Mello, P. Naphthenic Acids
724 Aggregation: The Role of Salinity. *Computation* **2022**, **10** (10), 170.

- 725 (51) Tshepelevitsh, S.; Hernits, K.; Leito, I. Prediction of Partition and Distribution Coefficients
726 in Various Solvent Pairs with COSMO-RS. *J Comput Aided Mol Des* **2018**, *32* (6), 711–722.
- 727 (52) Sonia Losada-Barreiro, Fátima Paiva-Martins, Carlos Bravo-Díaz. Partitioning of
728 Antioxidants in Edible Oil–Water Binary Systems and in Oil-in-Water
729 Emulsions. *Antioxidants* **2023**, *12* (4), 828. <https://doi.org/10.3390/antiox12040828>
- 730 (53) Escher BI, Abagyan R, Embry M, Klüver N, Redman AD, Zarfl C, Parkerton TF.
731 Recommendations for Improving Methods and Models for Aquatic Hazard Assessment of
732 Ionizable Organic Chemicals. *Environ Toxicol Chem.* 2020 Feb;39(2):269-286. doi:
733 10.1002/etc.4602.
- 734 (54) Hansima MACK, Zvomuya F, Amarakoon I. Fate of veterinary antimicrobials in Canadian
735 prairie soils - A critical review. *Science of The Total Environment* **2023**, *892*,
736 164387. <https://doi.org/10.1016/j.scitotenv.2023.164387>
- 737 (55) Tsantili-Kakoulidou, A.; Panderi, I.; Csizmadia, F.; Darvas, F. Prediction of Distribution
738 Coefficient from Structure. 2. Validation of Prolog D, an Expert System. *American*
739 *Pharmaceutical Association* **1997**, *86* (10), 1173–1179.
- 740 (56) Zamora, W. J.; Bertsch, E.; Suñer, S.; Pinheiro, S. Experimental n-Octanol/Water
741 Partition/Distribution Coefficients Database for Small Molecules. **2023**.
742 <https://doi.org/10.5281/ZENODO.7956685>.
- 743 (57) Kim, S.; Chen, J.; Cheng, T.; Gindulyte, A.; He, J.; He, S.; Li, Q.; Shoemaker, B. A.;
744 Thiessen, P. A.; Yu, B.; Zaslavsky, L.; Zhang, J.; Bolton, E. E. PubChem in 2021: New Data
745 Content and Improved Web Interfaces. *Nucleic Acids Res* **2021**, *49* (D1), D1388–D1395.
746 <https://doi.org/10.1093/NAR/GKAA971>
- 747 (58) Mishra, B.; Sankar, C.; Mishra, M. Polymer Based Solutions of Bupranolol Hydrochloride
748 for Intranasal Systemic Delivery. *J Drug Target* **2011**, *19* (3), 204–211.
749 <https://doi.org/10.3109/1061186X.2010.492520>.
- 750 (59) Bello, M. L.; Junior, A. M.; Freitas, C. A.; Moreira, M. L. A.; da Costa, J. P.; de Souza, M.
751 A.; Santos, B. A. M. C.; de Sousa, V. P.; Castro, H. C.; Rodrigues, C. R.; Cabral, L. M.
752 Development of Novel Montmorillonite-Based Sustained Release System for Oral Bromopride
753 Delivery. *European Journal of Pharmaceutical Sciences* **2022**, *175*.
754 <https://doi.org/10.1016/j.ejps.2022.106222>.
- 755 (60) Bezençon, J.; Wittwer, M. B.; Cutting, B.; Smieško, M.; Wagner, B.; Kansy, M.; Ernst, B.
756 PKa Determination by ¹H NMR Spectroscopy - An Old Methodology Revisited. *J Pharm Biomed*
757 *Anal* **2014**, *93*, 147–155. <https://doi.org/10.1016/j.jpba.2013.12.014>.
- 758 (61) Voigt, W.; Mannhold, R.; Limberg, J.; Blaschke, G. Interactions of Antiarrhythmics with
759 Artificial Phospholipid Membranes. *J Pharm Sci* **1988**, *77* (12), 1018–1020.
- 760 (62) Roseman, T. J.; Yalkowsky, S. H. Physicochemical Properties of Prostaglandin F_{2a}
761 (Tromethamine Salt): Solubility Behavior, Surface Properties, and Ionization Constants. *J Pharm*
762 *Sci* **1973**, *62* (10), 1680–1685.
- 763 (63) Shalaeva, M.; Kenseth, J.; Lombardo, F.; Bastin, A. Measurement of Dissociation
764 Constants (PK_a Values) of Organic Compounds by Multiplexed Capillary Electrophoresis Using

- 765 Aqueous and Cosolvent Buffers. *J Pharm Sci* **2008**, *97* (7), 2581–2606.
766 <https://doi.org/10.1002/jps.21287>.
- 767 (64) Qiang, Z.; Adams, C. Potentiometric Determination of Acid Dissociation Constants (PK a)
768 for Human and Veterinary Antibiotics. *Water Res* **2004**, *38* (12), 2874–2890.
769 <https://doi.org/10.1016/j.watres.2004.03.017>.
- 770 (65) de Almeida Drumond dos Santos, T.; Oliveira da Costa, D.; Silva da Rocha Pita, S.; Silva
771 Semaan, F. Potentiometric and Conductimetric Studies of Chemical Equilibria for Pyridoxine
772 Hydrochloride in Aqueous Solutions: Simple Experimental Determination of PKa Values and
773 Analytical Applications to Pharmaceutical Analysis. *Ecl. Quím* **2010**, *35* (4), 81–86.
- 774 (66) Morimoto, K.; Nagayasu, A.; Fukanoki, S.; Morisaka, K.; Hyon, S.-H.; Ikada, Y.
775 Evaluation of Polyvinyl Alcohol Hydrogel as Sustained-Release Vehicle for Transdermal System
776 of Bunitrolol-HCl-1. *Drug Dev Ind Pharm* **1990**, *16* (1), 13–29.
- 777 (67) Loftsson, T.; Thorisdóttir, S.; Fridriksdóttir, H.; Stefánsson, E. Enalaprilat and Enalapril
778 Maleate Eyedrops Lower Intraocular Pressure in Rabbits. *Acta Ophthalmol* **2010**, *88* (3), 337–341.
779 <https://doi.org/10.1111/j.1755-3768.2008.01495.x>.
- 780 (68) Mannhold, R.; Dross, K. P.; Frekker, R.; der Steen, van. Drug Lipophilicity in QSAR
781 Practice: I. A Comparison of Experimental with Calculative Approaches. *Quant. Struct. Act. Relat.*
782 **1990**, *9*, 21–28.
- 783 (69) Loftsson, T.; Vogensen, S. B.; Desbos, C.; Jansook, P. Carvedilol: Solubilization and
784 Cyclodextrin Complexation: A Technical Note. *AAPS PharmSciTech* **2008**, *9* (2), 425–430.
785 <https://doi.org/10.1208/s12249-008-9055-7>.
- 786 (70) Kuntworbe, N.; Alany, R. G.; Brimble, M.; Al-Kassas, R. Determination of PKa and
787 Forced Degradation of the Indoloquinoline Antimalarial Compound Cryptolepine Hydrochloride.
788 *Pharm Dev Technol* **2013**, *18* (4), 866–876. <https://doi.org/10.3109/10837450.2012.668554>.
- 789 (71) Islam, M. S.; Narurkar, M. M. Solubility, Stability and Ionization Behaviour of Famotidine.
790 *Journal of Pharmacy and Pharmacology* **1993**, *45* (8), 682–686. <https://doi.org/10.1111/j.2042-7158.1993.tb07088.x>.
- 792 (72) Deng, Y.; Li, B.; Yu, K.; Zhang, T. Biotransformation and Adsorption of Pharmaceutical
793 and Personal Care Products by Activated Sludge after Correcting Matrix Effects. *Science of the*
794 *Total Environment* **2016**, *544*, 980–986. <https://doi.org/10.1016/j.scitotenv.2015.12.010>.
- 795 (73) Franke, U.; Munk, A.; Wiese, M. Ionization Constants and Distribution Coefficients of
796 Phenothiazines and Calcium Channel Antagonists Determined by a PH-Metric Method and
797 Correlation with Calculated Partition Coefficients. *J Pharm Sci* **1999**, *88* (1), 89–95.
798 <https://doi.org/10.1021/js980206m>.
- 799 (74) Avdeef, A.; Box, K. J.; Comer, J. E. A.; Hibbert, C.; Tam, K. Y. PH-Metric LogP 10.
800 Determination of Liposomal Membrane-Water Partition Coefficient of Ionizable Drugs. *Pharm*
801 *Res* **1998**, *15* (2), 209–215.
- 802 (75) Thanacoody, R. H. K. Thioridazine: The Good and the Bad. *Recent Pat Antiinfect Drug*
803 *Discov* **2011**, *6*, 92–98.

- 804 (76) Martínez, V.; Maguregui, M. I.; Jiménez, R. M.; Alonso, R. M. Determination of the PK a
805 Values of B-Blockers by Automated Potentiometric Titrations. *J Pharm Biomed Anal* **2000**, *23*,
806 459–468.
- 807 (77) Huerta, B.; Jakimska, A.; Gros, M.; Rodríguez-Mozaz, S.; Barceló, D. Analysis of Multi-
808 Class Pharmaceuticals in Fish Tissues by Ultra-High-Performance Liquid Chromatography
809 Tandem Mass Spectrometry. *J Chromatogr A* **2013**, *1288*, 63–72.
810 <https://doi.org/10.1016/j.chroma.2013.03.001>.
- 811 (78) Fini, A.; Fazio, G.; Feroci, G. Solubility and Solubilization Properties of Non-Steroidal
812 Anti-Inflammatory Drugs. *Int J Pharm* **1995**, *126* (1–2), 95–102. [https://doi.org/10.1016/0378-](https://doi.org/10.1016/0378-5173(95)04102-8)
813 [5173\(95\)04102-8](https://doi.org/10.1016/0378-5173(95)04102-8).
- 814 (79) Jacka, M. R. *Clarke's Isolation and Identification of Drugs*, 2nd ed.; Moffat, A. C.,
815 Jackson, J. V., Moss, M. S., Widdop, B., Greenfield, E. S., Eds.; Pharmaceutical Press: London,
816 2000.
- 817 (80) Nakamura, Y.; Yamamoto, H.; Sekizawa, J.; Kondo, T.; Hirai, N.; Tatarazako, N. The
818 Effects of PH on Fluoxetine in Japanese Medaka (*Oryzias Latipes*): Acute Toxicity in Fish Larvae
819 and Bioaccumulation in Juvenile Fish. *Chemosphere* **2008**, *70* (5), 865–873.
820 <https://doi.org/10.1016/j.chemosphere.2007.06.089>.
- 821 (81) Schröder, W.; Andersson, J. T. Fast and Direct Method for Measuring 1-Octanol-Water
822 Partition Coefficients Exemplified for Six Local Anesthetics. *J Pharm Sci* **2001**, *90* (12), 1948–
823 1954. <https://doi.org/10.1002/JPS.1145>.
- 824 (82) Avdeef, A. *Sirius Technical Application Notes (STAN)*; Sirius Analytical Instruments Ltd.:
825 Forest Row, UK, 1994; Vol. 1.
- 826 (83) Avdeef, A.; Box, K. J.; Comer, J. E. A.; Hibbert, C.; Tam, K. Y. PH-Metric LogP 10.
827 Determination of Liposomal Membrane-Water Partition Coefficients of Ionizable Drugs. *Pharm*
828 *Res* **1998**, *15* (2), 209–215. <https://doi.org/10.1023/A:1011954332221/METRICS>.
- 829 (84) Caron, G.; Steyaert, G.; Pagliara, A.; Âde, F.; Reymond, Â.; Crivori, P.; Gaillard, P.;
830 Carrupt, P.-A.; Avdeef, A.; Comer, J.; Box, K. J.; Girault, H. H.; Testa, B. Structure-Lipophilicity
831 Relationships of Neutral and Protonated b-Blockers Intra- and Intermolecular Effects in Isotropic
832 Solvent Systems. [https://doi.org/10.1002/\(SICI\)1522-2675\(19990804\)82:8](https://doi.org/10.1002/(SICI)1522-2675(19990804)82:8).
- 833 (85) Avdeef, A. *Sirius Technical Application Notes (STAN)*; Sirius Analytical Instruments Ltd.:
834 Forest Row, UK, 1995; Vol. 2.
- 835 (86) Lombardo, F.; Shalaeva, M. Y.; Tupper, K. A.; Gao, F.; Abraham, M. H. ElogP(Oct): A
836 Tool for Lipophilicity Determination in Drug Discovery. *J Med Chem* **2000**, *43* (15), 2922–2928.
837 <https://doi.org/10.1021/JM0000822/ASSET/IMAGES/MEDIUM/JM0000822E00013.GIF>.
- 838 (87) Winiwarter, S.; Bonham, N. M.; Ax, F.; Hallberg, A.; Lennernäs, H.; Karlén, A.
839 Correlation of Human Jejunal Permeability (in Vivo) of Drugs with Experimentally and
840 Theoretically Derived Parameters. A Multivariate Data Analysis Approach. *J Med Chem* **1998**, *41*
841 (25), 4939–4949. https://doi.org/10.1021/JM9810102/SUPPL_FILE/JM9810102_S.PDF.
- 842 (88) Slater, B.; McCormack, A.; Avdeef, A.; Comer, J. E. A. PH-Metric Log P. 4. Comparison
843 of Partition Coefficients Determined by HPLC and Potentiometric Methods to Literature Values.
844 *J Pharm Sci* **1994**, *83* (9), 1280–1283. <https://doi.org/10.1002/JPS.2600830918>.

845 (89) Luger, P.; Daneck, K.; Engel, W.; Trummnitz, G.; Wagner, K. Structure and
846 Physicochemical Properties of Meloxicam, a New NSAID. *European Journal of Pharmaceutical*
847 *Sciences* **1996**, *4* (3), 175–187. [https://doi.org/10.1016/0928-0987\(95\)00046-1](https://doi.org/10.1016/0928-0987(95)00046-1).

848 (90) Takács-Novák, K.; Józán, M.; Hermeicz, I.; Szász, G. Lipophilicity of Antibacterial
849 Fluoroquinolones. *Int J Pharm* **1992**, *79* (1–3), 89–96. [https://doi.org/10.1016/0378-](https://doi.org/10.1016/0378-5173(92)90099-N)
850 [5173\(92\)90099-N](https://doi.org/10.1016/0378-5173(92)90099-N).

851 (91) Carda-Broch, S.; Berthod, A. PH Dependence of the Hydrophobicity of β -Blocker Amine
852 Compounds Measured by Counter-Current Chromatography. *J Chromatogr A* **2003**, *995* (1–2),
853 55–66. [https://doi.org/10.1016/S0021-9673\(03\)00534-X](https://doi.org/10.1016/S0021-9673(03)00534-X).

854 (92) Scott, D. Estimation of Distribution Coefficients from the Partition Coefficient and PKa.
855 *Pharmaceutical Technology* **2002**, *26* (11).

856 (93) Gulaboski, R.; Borges, F.; Pereira, C. M.; Natália, M.; Cordeiro, D. S.; Garrido, J.; Silva,
857 A. F. Voltammetric Insights in the Transfer of Ionizable Drugs Across Biomimetic Membranes-
858 Recent Achievements. *Comb Chem High Throughput Screen* **2007**, *10*, 514–526.

859 (94) *pKa Plugin | Chemaxon Docs*. <https://docs.chemaxon.com/display/docs/pka-plugin.md>
860 (accessed 2023-01-24).

861 (95) *LogP and logD calculations | Chemaxon Docs*.
862 <https://docs.chemaxon.com/display/docs/logp-and-logd-calculations.md> (accessed 2023-01-24).

863 (96) <https://cran.r-project.org/web/packages/Metrics/Metrics.pdf>

864 (97) E. L. Willighagen, J. W. Mayfield, J. Alvarsson, A. Berg, L. Carlsson, N. Jeliaskova, S.
865 Kuhn, T. Pluskal, M. Rojas-Chertó, O. Spjuth, G. Torrance, C. T. Evelol, R. Guha and G.
866 Steinbeck, *J. Cheminform.*, **2017**, *9*, 1–19.

867 (98) Ghiandoni, G. M.; Caldeweyher, E. Fast Calculation of Hydrogen-Bond Strengths and Free
868 Energy of Hydration of Small Molecules. *Scientific Reports* **2023**, *13*:1 **2023**, *13* (1), 1–11.
869 <https://doi.org/10.1038/S41598-023-30089-X>.

870 (99) Welch, B. L. The Generalization Of ‘Student’s’ Problem When Several Different
871 Population Variances Are Involved. *Biometrika* **1947**, *34* (1–2), 28–35.
872 <https://doi.org/10.1093/BIOMET/34.1-2.28>.

873 (100) Ruxton, G. D. The Unequal Variance T-Test Is an Underused Alternative to Student’s t-
874 Test and the Mann–Whitney U Test. *Behavioral Ecology* **2006**, *17* (4), 688–690.
875 <https://doi.org/10.1093/BEHECO/ARK016>.

876 (101) Kuhn, M.; Johnson, K. *Feature Engineering and Selection: A Practical Approach for*
877 *Predictive Models*, 1st ed.; Taylor and Francis Group: London, 2019.

878 (102) <https://cran.r-project.org/web/packages/caret/>

879 (103) Liu, S.; Cao, C.; Li, Z. Approach to Estimation and Prediction for Normal Boiling Point
880 (NBP) of Alkanes Based on a Novel Molecular Distance-Edge (MDE) Vector, γ . *J Chem Inf*
881 *Comput Sci* **1998**, *38* (3), 387–394.

882 (104) Ghose, A. K.; Crippen, G. M. Atomic Physicochemical Parameters for Three-Dimensional
883 Structure-Directed Quantitative Structure-Activity Relationships I. Partition Coefficients as a

884 Measure of Hydrophobicity. *J Comput Chem* **1986**, 7 (4), 565–577.
885 <https://doi.org/10.1002/JCC.540070419>.

886 (105) Šegatin, N.; Klofutar, C. Thermodynamics of the Solubility of Water in 1-Hexanol, 1-
887 Octanol, 1-Decanol, and Cyclohexanol. *Monatsh Chem* **2004**, 135 (3), 241–248.

888 (106) Hall, L. H.; Kier, L. B. Electrotological State Indices for Atom Types: A Novel
889 Combination of Electronic, Topological, and Valence State Information. *J Chem Inf Comput Sci*
890 **1995**, 35 (6), 1039–1045.

891 (107) <https://cran.r-project.org/web/packages/bestglm/bestglm.pdf>

892 (108) Furnival, G. M.; Wilson, R. W. Regressions by Leaps and Bounds. *Technometrics* **1974**,
893 16 (4), 499–511. <https://doi.org/10.1080/00401706.1974.10489231>.

894 (109) Burger, S. *Introduction to Machine Learning with R: Rigorous Mathematical Modeling*,
895 1st ed.; O'Reilly: United States of America, 2018.

896 (110) Fauchere, J. L. & Pliska, V. Hydrophobic Parameters π of Amino Acid Side Chains
897 from the Partitioning of *N*-Acetyl-*L*-Amino Acid Amides. *Eur. J. Med. Chem.* **1983**, **18**, 369–
898 375.

899 (111) Breiman, L. Random Forests. *Mach Learn* **2001**, 45 (1), 5–32.
900 <https://doi.org/10.1023/A:1010933404324/METRICS>.

901 (112) <https://cran.r-project.org/web/packages/randomForest/index.html>.

902 (113) Cortes, C.; Vapnik, V.; Saitta, L. Support-Vector Networks. *Machine Learning 1995* 20:3
903 **1995**, 20 (3), 273–297. <https://doi.org/10.1007/BF00994018>.

904 (114) Boser, B. E.; Guyon, I. M.; Vapnik, V. N. Training Algorithm for Optimal Margin
905 Classifiers. *Proceedings of the Fifth Annual ACM Workshop on Computational Learning Theory*
906 **1992**, 144–152. <https://doi.org/10.1145/130385.130401>.

907 (115) <https://cran.r-project.org/web/packages/e1071/index.html>

908 (116) Viayna A, Antermite SG, de Candia M, Altomare CD, Luque FJ. Interplay between
909 Ionization and Tautomerism in Bioactive β -Enamino Ester-Containing Cyclic Compounds: Study
910 of Annulated 1,2,3,6-Tetrahydroazocine Derivatives. *J Phys Chem B.* **2020** Jan 9;124(1):28-37.
911 doi: 10.1021/acs.jpcc.9b08904.

912 (117) Tielker N, Güssregen S, Kast SM. SAMPL7 physical property prediction from EC-RISM
913 theory. *J Comput Aided Mol Des.* **2021** Aug;35(8):933-941. doi: 10.1007/s10822-021-00410-9.

914 (118) Viayna A, Pinheiro S, Curutchet C, Luque FJ, Zamora WJ. Prediction of n-octanol/water
915 partition coefficients and acidity constants (pK_a) in the SAMPL7 blind challenge with the
916 IEFPCM-MST model. *J Comput Aided Mol Des.* **2021** Jul;35(7):803-811.

917 (119) Rodriguez SA, Tran JV, Sabatino SJ, Paluch AS. Predicting octanol/water partition
918 coefficients and pK_a for the SAMPL7 challenge using the SM12, SM8 and SMD solvation models.
919 *J Comput Aided Mol Des.* **2022** Sep;36(9):687-705. doi: 10.1007/s10822-022-00474-1.

920 (120) Wu J, Kang Y, Pan P, Hou T. Machine learning methods for pK_a prediction of small
921 molecules: Advances and challenges. *Drug Discov Today.* **2022** Dec;27(12):103372. doi:
922 10.1016/j.drudis.2022.103372.

923 (121) Johnston RC, Yao K, Kaplan Z, Chelliah M, Leswing K, Seekins S, Watts S, Calkins D,
924 Chief Elk J, Jerome SV, Repasky MP, Shelley JC. Epik: pK_a and Protonation State Prediction
925 through Machine Learning. *J Chem Theory Comput.* **2023**, 25;19(8):2380-2388. doi:
926 10.1021/acs.jctc.3c00044.

927 (122) Reymond, F.; Chopineaux-Courtois, V.; Steyaert, G.; Bouchard, G.; Carrupt, P. A.; Testa,
928 B.; Girault, H. H. Ionic Partition Diagrams of Ionisable Drugs: PH-Lipophilicity Profiles, Transfer
929 Mechanisms and Charge Effects on Solvation. *Journal of Electroanalytical Chemistry* **1999**, 462
930 (2), 235–250. [https://doi.org/10.1016/S0022-0728\(98\)00418-5](https://doi.org/10.1016/S0022-0728(98)00418-5)

931 (123) Burden, F. R. Molecular Identification Number for Substructure Searches. *J Chem Inf*
932 *Comput Sci* **1989**, 29 (3), 225–227.

933 (124) Burden, F. R. A Chemically Intuitive Molecular Index Based on the Eigenvalues of a
934 Modified Adjacency Matrix. *Quant. Stmct-Act. Relat* **1997**, 16, 3–314.

935 (125) Pearlman, R. S.; Smith, K. M. Metric Validation and the Receptor-Relevant Subspace
936 Concept. *J Chem Inf Comput Sci* **1999**, 39 (1), 28–35.

937 (126) Petitjean, M. Applications of the Radius-Diameter Diagram to the Classification of
938 Topological and Geometrical Shapes of Chemical Compounds. *J Chem Inf Comput Sci* **1992**, 32
939 (4), 331–337.

940 (127) Ertl, P.; Rohde, B.; Selzer, P. Fast Calculation of Molecular Polar Surface Area as a Sum of
941 Fragment-Based Contributions and Its Application to the Prediction of Drug Transport Properties.
942 *J Med Chem* **2000**, 43 (20), 3714–3717.

943

944

945

946

947

948

949

950

951

952

953

954

955

956

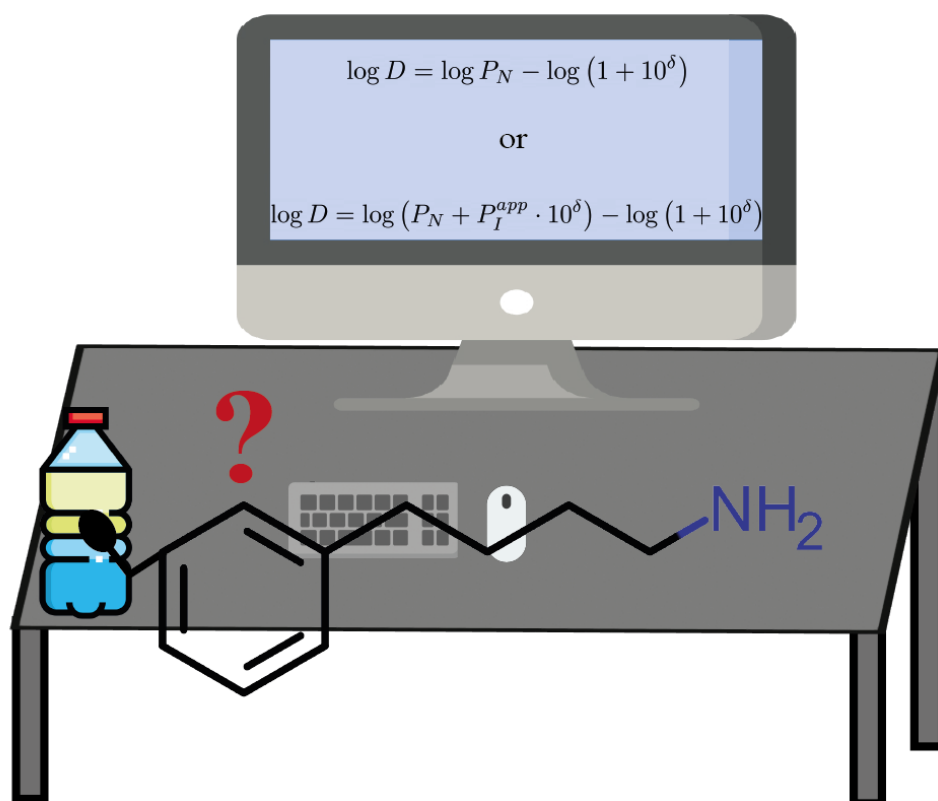
957

958

959 *TOC Graphics*

960

961



962