

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

3 Esteban Bertsch¹⁺, Sebastián 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 José, Costa Rica

10

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

12 ⁺These authors contributed equally to this work

13

14 **Keywords** Partition coefficient, Lipophilicity profiles, Machine learning, Chemoinformatics,
15 Drug 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 is
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 both the experimental lipophilicity
39 of small molecules and experimental lipophilicity-based applications and metrics such as lipophilic
40 efficiency (LipE), distribution of spiked drugs in milk products, and pH-dependent partition of
41 water contaminants in synthetic passive samples such as silicones. Our findings show that better
42 predictions are obtained by considering the apparent ionic partition, whereas ignoring its
43 contribution can lead to inadequate experimental simplifications and/or computational predictions.
44 In this context, we developed machine learning algorithms to determine the cases that P_I^{app} should
45 be considered. The results indicate that small, rigid, and unsaturated molecules with $\log P_N$ close
46 to zero, which present a significant proportion of ionic species in the aqueous phase, were better
47 modeled using the apparent ionic partition (P_I^{app}). In addition, we validated our findings using a
48 test and two external sets, which included small molecules and amino acid analogs, where the
49 logistic regressions, random forest classifications, and support vector machine models predicted
50 better formalism to determine the $\log D_{\text{pH}}$ for each molecule with high accuracies, sensitivities, and
51 specificities. Finally, our findings can serve as guidance to the scientific community working in
52 early-stage drug design, food, and environmental chemistry who deal with ionizable molecules, to
53 determine a priori which pH-dependent lipophilicity profile should be used in their research and
54 applications depending on the structure of a substance.

55

56

57

58 Introduction

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

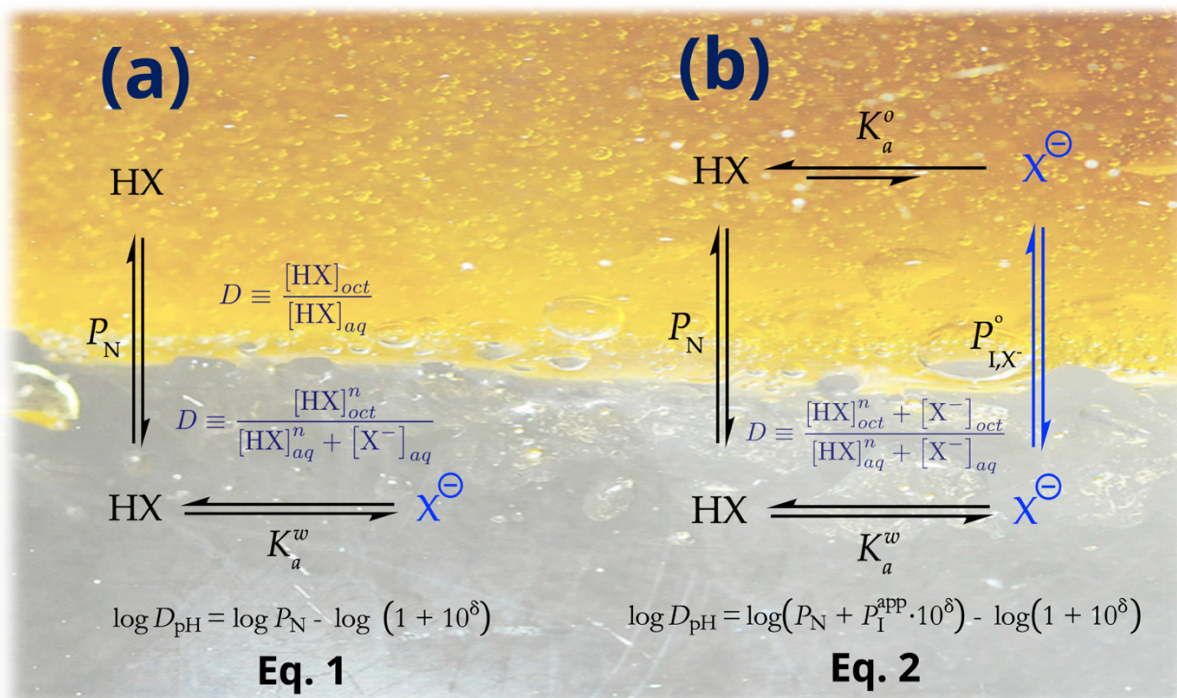
71 Since the pH of the solution directly affects the concentration of neutral and ionic species,
72 the equilibrium constant varies with pH, which also means that the lipophilicity of a compound is
73 dependent on it. The partition coefficient as a function of pH is often called distribution coefficient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

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

88
$$\log D_{\text{pH}} = \log P_{\text{N}} - \log (1 + 10^{\delta}) \quad [1]$$

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

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



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

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

104
$$\log D_{\text{pH}} = \log(P_{\text{N}} + P_{\text{I}}^{\text{app}} \cdot 10^{\delta}) - \log(1 + 10^{\delta})$$
 [2]

105

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

126 Despite the lack of a consensus formalism to model $\log D_{\text{pH}}$ as a function of $P_{\text{I}}^{\text{app}}$, and
127 considering that different theoretical approaches have shown similar trends^{14,21,22}, Equation 2 has
128 been successfully used for modeling the lipophilicity of ionized compounds in many areas of basic
129 and applied sciences. For instance, to study the 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 has been 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,
139 whereas Equation 2 predicts a $\log D_{\text{pH}}$ value closer to the experimental values. This study showed
140 that Equation 2 provided a more exact lipophilicity profile over a wider pH range than Equation
141 1. However, no systematized study has been performed to evaluate the importance of considering
142 the ionic partition on the $\log D_{\text{pH}}$ prediction for large sets of small drug-like molecules at various
143 pH values, although it has been reported that much of the poor performance of some models on
144 blind challenges has been due to the simplification of ignoring the ionic species partition.²⁷

145 In this study, our aim is to evaluate the impact of considering the $P_{\text{I}}^{\text{app}}$ in determining pH-
146 dependent lipophilicity profiles of small molecules. We also aim to provide guidance to the
147 scientific community working in early-stage drug design, food, and environmental chemistry,
148 specifically those dealing with ionizable molecules. Our goal is to help researchers determine a
149 priori which pH-dependent lipophilicity profile should be used based solely on structural features
150 of the substance of interest. To this end, we collected the experimental values of $\log P_{\text{N}}$, $\text{p}K_{\text{a}}$, and
151 $\log P_{\text{I}}^{\text{app}}$ of different compounds at various pH values as well as experimental data of lipophilicity-
152 based applications and metrics such as lipophilic efficiency (LipE), distribution of spiked drugs in
153 milk products and pH-dependent partition in passive samples, which were used to compute $\log D_{\text{pH}}$
154 with Equations 1 and 2. The predictions using both equations were then used to compare their
155 performances using statistical parameters. Finally, logistic regression (**LR**), random forest
156 classification (**RFC**), and support vector machine (**SVM**) models were developed to define from
157 the molecular structure which formalism is recommended for modeling pH-dependent lipophilicity
158 profiles.

159

160 Methodology

161 Data collection and classification

162 We critically compiled the experimental values of $\log P_N$, pK_a , $\log P_I^{\text{app}}$, and $\log D_{\text{pH}}$ of 225
163 entries based on earlier literature reports (database available in reference 33).^{29,55,56} Refs. 29 and
164 55 were chosen based on the wide selection of experimental data for $\log P_N$, $\log D_{\text{pH}}$, and $\log P_I$
165 values and because they encompass the desired chemical space of small molecules for our
166 modeling. SMILES codes were collected from publicly available data in PubChem.⁵⁷ The
167 experimental pK_a values were also obtained from PubChem, but they were corroborated by
168 reviewing their values in primary literature reports.^{38,57-80} The experimental technique of $\log P_N$,
169 $\log D_{\text{pH}}$, and $\log P_I^{\text{app}}$ measurements for each entry were thoroughly revised and added to the
170 database.^{74,81-90} Ref 55 provided experimental $\log D_{\text{pH}}$ values of molecules in diverse pH ranges.
171 The $\log P_I$ values were obtained from the $\log D_{\text{pH}}$ at the most extreme measured pH, in which the
172 molecule would be mostly (above 95 %) in its ionized state. The $\log P_I^{\text{app}}$ values for molecules that
173 were not measured under ionizable pH conditions were obtained from external sources.^{38,74,91,92}
174 The molecules were classified as acids or bases based on their functional groups and pK_a values.
175 Zwitterionic compounds were found by evaluating the difference between acidic and basic pK_a in
176 conjunction with ChemAxon's calculator of protonated species distribution in function of pH.⁹³
177 Zwitterionic and amphoteric species were also classified as acidic or basic based on the behavior
178 of their lipophilicity profiles, which were evaluated using the ChemAxon partitioning calculator.⁹⁴

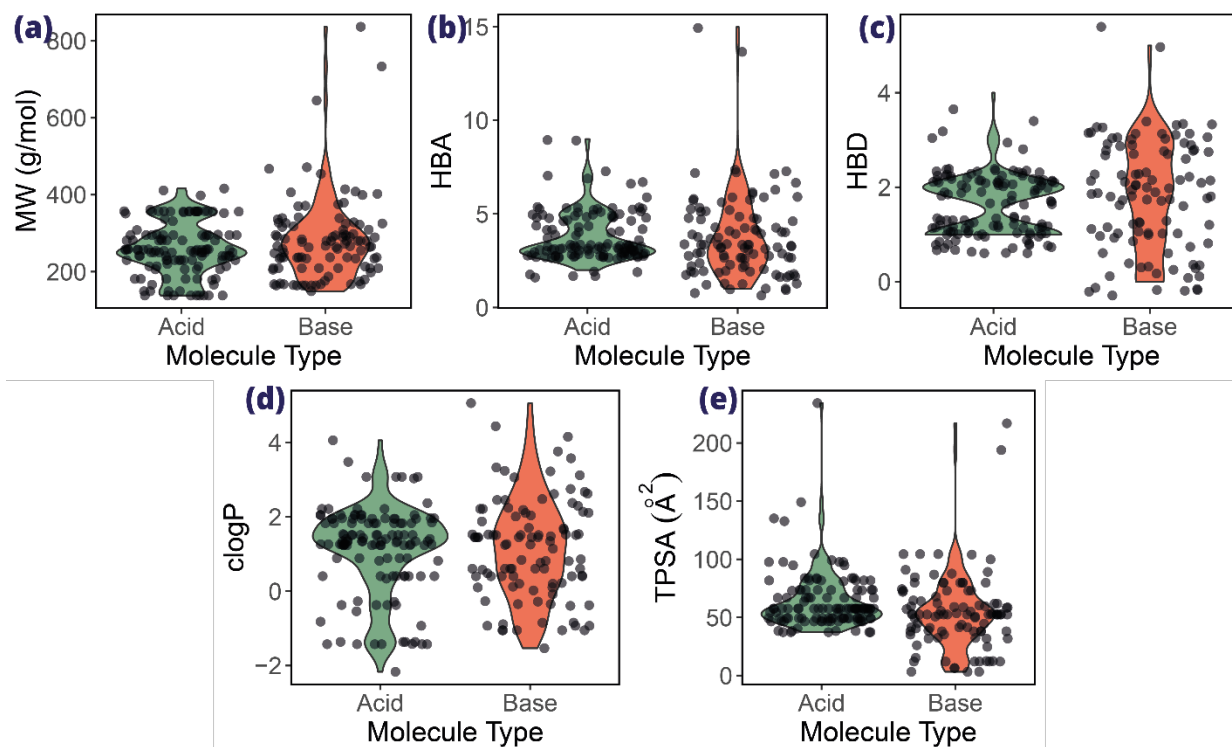
179 Figure 2 shows the distribution of the molecules along several descriptors of their chemical
180 space. Most compounds can be considered small molecules because they tend to have small
181 molecular weights ($< 400 \text{ gmol}^{-1}$) and topological polar surface areas ($< 100 \text{ \AA}^2$). Our database
182 consists mostly of lipophilic species since the $\text{clog}P$ values are mostly positive, which coincides
183 with the low polarity of our molecules, demonstrated by the tendency of low counts of hydrogen
184 bond donors (< 3) and acceptors (< 7).

185

186

187

188



189

190 Figure 2. Distribution of molecular properties in the database⁵⁶ by (a) Molecular weight (MW), (b)
 191 hydrogen bond acceptors, (c) hydrogen bond donors, (d) calculated $\log P_N$ (obtained with Alogp)⁹⁵,
 192 and (d) topological polar surface area. These descriptors were calculated using the ‘*RCDK*’
 193 package in R.

194

195 Performance of pH-dependent lipophilicity profiles

196

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

$$200 \quad d_1 = \left| \log D_{\text{Eq.1}} - \log D_{\text{exp}} \right| \quad [3]$$

$$201 \quad d_2 = \left| \log D_{\text{Eq.2}} - \log D_{\text{exp}} \right| \quad [4]$$

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

203 The performance of the two formalisms was tested by performing a linear regression of
204 $\log D_{\text{Eq.1}}$ and $\log D_{\text{Eq.2}}$ on their experimental values. The root-mean-squared error (RMSE), mean
205 absolute error (MAE), mean squared error (MSE), and Pearson's correlation coefficient squared
206 (R^2) were calculated with the 'Metrics' package in R.⁹⁶ We also tested the performance of each
207 formalism on each individual molecule using the descriptor d_3 (Eq. 5). When d_3 yielded a value
208 greater than zero, Eq. 2 fits better than Eq. 1. Accordingly, we created a binomial conditional based
209 on the values of d_3 , where Eq. 2 should be used when d_3 is greater than 0.2 (see Results and
210 Discussion); otherwise, both equations are considered equivalent, which can be interpreted as Eq.
211 1 provides better modeling, owing to its simplicity.

$$212 \quad d_3 = d_1 - d_2 \quad [5]$$

213

214 **Experimental data of lipophilicity-based applications and metrics used in medicinal, food,** 215 **and environmental chemistry.**

216 We also investigated the impact of the apparent ionic partitioning contribution to
217 lipophilicity-based parameters commonly used in the fields of food, medicinal, and environmental
218 chemistry. Two tests were conducted for food applications. First, we evaluated Eqs. 1 and 2 to
219 reproduce the experimental $\log D_{4.5}$ for the partition of bioantioxidants in a oil/water system.⁹⁷
220 Secondly, we collected data on the distribution of spiked drugs in milk products using the pH, pK_a ,
221 $\log P_N$ and $\log D_{6.8}$ reported in the original work.^{98,99} However, ionic partition was obtained from
222 ChemAxon, except for the oxytetracycline (OTET), for which the experimental $\log P_1^{\text{app}}$ was found
223 in the literature¹⁰⁰ and used to measure the $\log D_{6.8}$ using Eq. 1 and Eq. 2.

224 In addition, as an environmentally relevant application, we obtained experimental pH-
225 dependent distribution data for a series of ionizable compounds on a passive sampler
226 polydimethylsiloxane (PDMS) and water. For this task, monoprotic acids and bases were searched
227 for within the 514 compounds in the article. The experimental pK_a values, $\log D_{\text{PDMS/w}}$, at several
228 pH values including extreme ranges (from which we were able to obtain $\log P_N$ and $\log P_1^{\text{app}}$) were
229 provided by the article.¹⁰¹ Therefore, predictions of the distribution coefficients in the PDMS/water
230 system to pH = 7.4 using Eq. 1 and Eq. 2. were calculated and compared with those reported in
231 experimental work.

232 Finally, we explored the influence of Eq 1. and Eq. 2 to predict a relevant metric used in
233 medicinal chemistry lead optimization affairs, the lipophilic efficiency (LipE):¹⁰²

$$234 \quad \text{LipE} = pAct - \log D_{pH}, \quad [6]$$

235 where $\log D_{pH}$ stands for the distribution coefficient and $pAct$ represents the negative
236 logarithm of biological activity, that is, half maximal inhibitory concentration (IC_{50} , mol/L),
237 inhibitory constant (K_i , mol/L), or binding energy constant (K_b). Here, we searched the literature
238 for ionizable monoacidic or monobasic drug-like molecules with both experimentally determined
239 $\log D_{pH}$ measurements and biological activities.^{103–114} In some cases, the experimental $\log P_N$, pK_a ,
240 and $\log P_I^{app}$ were reported, but otherwise they were determined using ChemAxon. The LipE was
241 then simulated using Eq 1. and Eq. 2 and compared to their experimental values.

242

243 **Machine Learning models to classify the molecules according to the best fit to pH-dependent** 244 **lipophilicity profiles**

245 Topological and constitutional descriptors were calculated with the software ‘*rdck*’
246 package in R¹¹⁵ while experimental measurements (i.e., $\log P_N$, pK_a , and pH) were added from our
247 dataset. We also added the free energies of hydration and hydrogen bond strengths computed using
248 the new open-source tool ‘*Jazzy*’.¹¹⁶ The H-bond donor and acceptor strengths were obtained by
249 calculating the partial charges of the hydrogen atoms and atoms with lone electron pairs,
250 respectively, along with corrective terms. The free energy of hydration was calculated using the
251 sum of the polar, apolar, and interaction terms. The polar term was derived from the previously
252 calculated H-bond donor and acceptor strengths. The apolar terms consist of the sum of the
253 weighted contributions of the topological surface area, number of rings, and p-orbital counts in the
254 sp and sp² atoms. The interaction term consists of a weighted sum of the amount of neighboring
255 H-bond acceptor groups each atom has in a molecule.¹¹⁶

256 We eliminated intercorrelated properties so that no descriptor had a correlation value of r^2
257 > 0.6 (Figure S1 and S2). After this filtration step, two different feature selection methods were
258 tested to choose the best descriptors for the Machine Learning models. First, we performed
259 Welch’s *t*-test (**WTT**), which evaluates the statistical difference between the means of two
260 populations that have unequal variances and sample sizes.^{117,118} The algorithm calculates the mean

261 of both groups from the binomial conditional for each descriptor. These values were evaluated
262 using Equation 7.

$$263 \quad t = \frac{\Delta\mu}{\delta_{\Delta\bar{x}}} \quad [7]$$

264
265 where t stands for the statistic t in Welch's t-test, and $\Delta\mu$ represents the mean difference between
266 data samples from each population (Eq. 1 or Eq. 2 better fits), and the uncertainty value of both
267 groups, which was calculated using the standard deviation of both population samples (Eq.8):

$$268 \quad \delta_{\Delta\bar{x}} = \sqrt{\left(\frac{s_1}{\sqrt{N_1}}\right)^2 + \left(\frac{s_2}{\sqrt{N_2}}\right)^2} \quad [8]$$

269 WTT was performed for each descriptor using R, and the p -value was extracted. Features that did
270 not show statistical significance between means ($p > 0.05$) were eliminated. Second, recursive
271 feature elimination (**RFE**) was performed. This iterative feature selection method builds a
272 predictive model using the entire set of descriptors and calculates its importance score (Figure S3).
273 The least important descriptors were removed, and the model was reiterated to achieve maximum
274 performance.¹¹⁹ This RFE algorithm was programmed using the ‘*caret*’ package in R¹²⁰ and tuned
275 via a 5-time repeated k -fold cross-validation ($k = 10$). Table 1 shows the descriptors selected using
276 the WTT feature selection method for acids and bases, along with their definitions and target
277 molecules. Table S1 lists the descriptors selected using the RFE method.

278

279 **Table 1.** List of the most influential structural descriptors^{95,116,121,122} used for the Machine Learning
 280 classification models, their target molecules, and the divergence between the two populations from
 281 our dataset were determined using the WTT feature selection method by separating the populations
 282 with the 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 a 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 |

283

284 **Logistic Regression Classification**

285 A logistic regression (LR) is a simple classification statistical model that provides a binary
286 response to the distribution of the input data among a specific descriptor. The simplest regressions
287 fit the distributions of data to a sigmoidal function, where the input values are given a probability
288 value, which is then classified into one of the two classes based on a cut-off value. We firstly
289 performed a feature selection process specific for logistic regressions by using the ‘*bestglm*’
290 package in R¹²³ which evaluates through n iterations, which combination of descriptors gives the
291 best fitted regression through the *leaps* algorithm.¹²⁴ This package evaluates the weight of each
292 descriptor by linearizing the sigmoidal function and giving a slope value and standard error for
293 each parameter like a multiple linear regression model (Equation 9).

$$294 \quad \ln\left(\frac{f(x)}{1-f(x)}\right) = \sum_{i=1}^n c_i x_i + b \quad [9]$$

295 The ‘*bestglm*’ package drops the parameters, where $c_i \rightarrow 0$. The algorithm iterates the
296 sigmoidal fit using Equation 8 n times until it finds the combination of descriptors in which the
297 parameters have the smallest standard error.¹²³ This feature selection process was performed
298 separately for acids and bases because the descriptors have different behaviors for each type of
299 molecule.

300 Figure 3 shows a flowchart of the modelling process. The dataset was divided into acids
301 (113 entries) and bases (100 entries). Zwitterions (12 entries) were not considered for the Machine
302 Learning predictions because of their small sample size and because further lipophilicity modeling
303 will be performed for these molecules (see Results and Discussion section). Acids and bases were
304 randomly and reproducibly sampled into the training and test sets at a ratio of 80:20. Multiple
305 logistic regressions were performed for the training sets based on previously collected descriptors.
306 Predictive models were programmed using the ‘*caret*’ package. Acids and bases were modeled
307 separately and labeled as **Models A** and **B**, respectively (see Figure 3). The test sets were evaluated
308 using both models. The performance of Models A and B was evaluated using confusion matrices
309 (see Table S2), which are widely used to evaluate classification models.¹²⁵ The confusion matrices
310 tabulate the number of true positives (**TP**), false positives (**FP**), true negatives (**TN**), and false
311 negative (**FN**) predictions, along with the sensitivity, specificity, and accuracy of the models.
312 Sensitivity determines the ability of the model to detect events of the positive class; that is, it
313 indicates the predictive performance of the molecules of the $\log D_{\text{Eq.2}}$ population (Equation 10). On

314 the other hand, the specificity indicates the performance of the model in detecting the negative
315 class, which in this case are the molecules of the $\log D_{\text{Eq.1}}$ population (Equation 11). The accuracy
316 indicates the overall performance in detecting false positives and false negatives (Equation 12).

317

$$318 \quad \text{Sensitivity} = \frac{TP}{TP + FN} \quad [10]$$

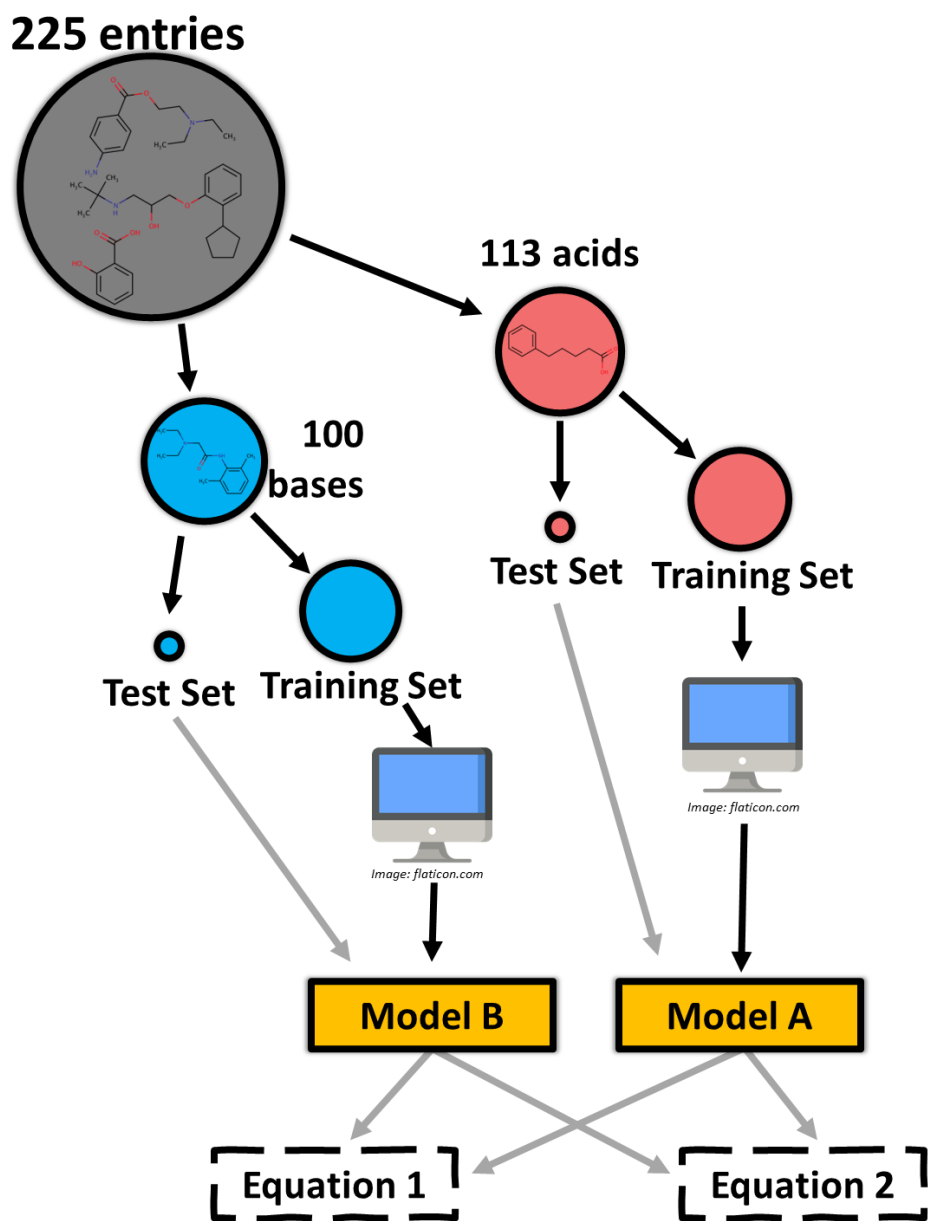
$$319 \quad \text{Specificity} = \frac{TN}{FP + TN} \quad [11]$$

$$320 \quad \text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad [12]$$

321

322 Models A and B were tested further using an external set. The experimental lipophilicity
323 measurements made by Disdier *et al.*⁴⁵ consisted of 69 data entries of small molecules with 38
324 acids, 16 bases, and 15 zwitterions, the latter being discarded for our analysis. To further check
325 the robustness of our models, a second external set of amino acid analogs was evaluated¹²⁶,
326 consisting of 8 entries of histidine (basic amino acid) and 10 entries of tyrosine (acidic amino acid).
327 Then, we evaluated the performance of Models A and B for this dataset using confusion matrices
328 (see Table S3-S4).

329



330

331 **Figure 3.** Graphical representation of the data classification and sampling of our dataset to create
 332 our predictive multiple logistic regression model using topological, constitutional, and
 333 experimental descriptors.

334

335

336 **Random Forest Classification**

337 Decision trees are a simple visual method for evaluating or classifying data, where each
338 node consists of a variable in the dataset. Each node leads to a leaf in which the desired output is
339 issued. A random forest is a combination of decision trees, which are randomly sampled, and the
340 nodes are randomly organized.¹²⁷ We split our dataset into training-, and test sets, as shown in
341 Figure 2. In this case, Models A and B consisted of random forest classifications (**RFC**) performed
342 with the ‘*randomForest*’ package in R.¹²⁸ Both models were previously refined using the *tuneRF*
343 function within the package, which chooses the optimal *mtry* variable. This value indicates the
344 number of features selected at each split in each decision tree, where *mtry* = 2 provided the best
345 prediction for both models (number of trees = 500, see Supporting Information Figure S4). The
346 importance of each descriptor in both models was evaluated through the mean decrease in the Gini
347 impurity index using the *MeanDecreaseGini* function (Figure S4).

348 The best lipophilicity profile fit for the acidic and basic tests and external sets was predicted with
349 Models A and B, respectively. The performance of each prediction was evaluated using confusion
350 matrices (see Tables S5-S7) and their respective sensitivity, specificity, and accuracy calculations
351 (Eqs. 10-12).

352

353 **Support Vector Machine Classification**

354 A Support Vector Machine (**SVM**) algorithm works by dividing training data into two
355 categories, either by linear or nonlinear classification; new data are then assigned to one of the two
356 classes. The model separates the data by finding a hyperplane that maximizes the gap between
357 categories. In the case of linear classification, the space is two-dimensional, making the hyperplane
358 a linear function.¹²⁹ When the data are not linearly separable, the algorithm performs the kernel
359 trick, which involves increasing the dimensions of the data space. This results in the hyperplane
360 being able to be another function in the original space, such as radial or polynomial, allowing to
361 classify the data in different ways.¹³⁰

362 We split our datasets in the same manner as the other classification models and set Models
363 A and B as support vector machines given by the ‘*e1071*’ package in R.¹³¹ We decided to compare
364 the performance of using a linear kernel (**SVML**) and a polynomial kernel (**SVMP**). Radial kernels
365 were not evaluated because our binary data do not follow a circular separation in the hyperplane;

366 therefore, our classifications do not provide an adequate fit. The hyperparameter selection for each
367 model was performed with the *trainControl* and *train* functions from the 'caret' package, which
368 executes a *k*-fold cross-validation (*k* = 10 was used), where different values of the parameters were
369 tested and selected, which resulted in the highest accuracy. The best hyperparameters were the
370 function's default parameters: *C* = 1 for SVMML and for SVMMP, *C* = 1, *degree* = 3, *gamma* = 1, and
371 *coef0* = 0. We calculated the accuracy, sensitivity, and specificity of each model using Eq. 9-11,
372 using the results from their respective confusion matrices (see Tables S8-S13). We then compared
373 the confusion matrices of the LR, RFC, SVMML, and SVMMP models to determine the one that
374 yielded the best results.

375

376 **Results and Discussion**

377

378 **Performance of pH-dependent lipophilicity profiles in predicting experimental distribution** 379 **coefficients**

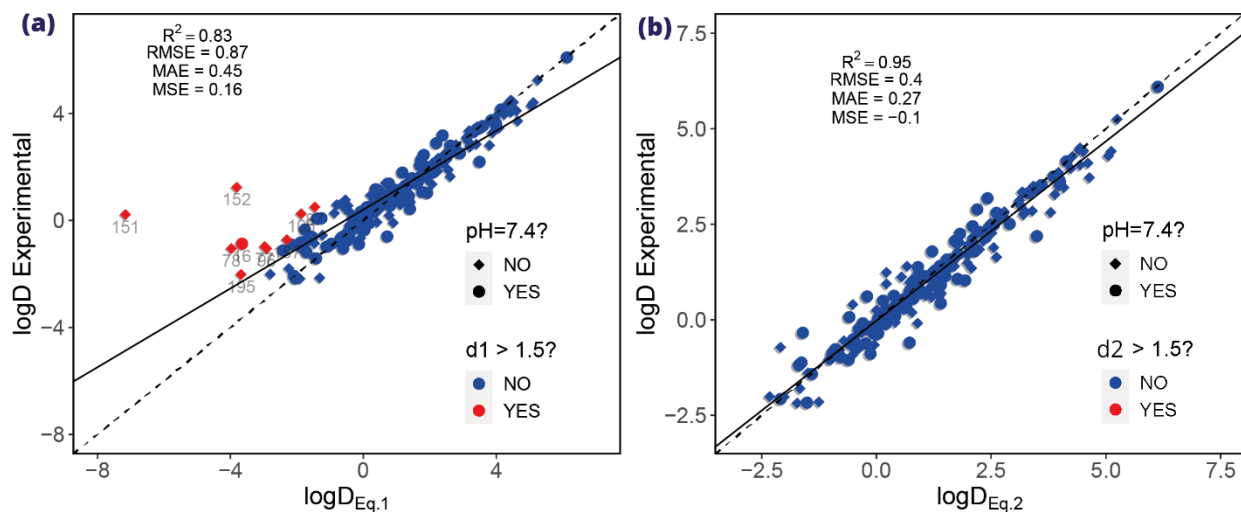
380

381 One of the main objectives of this study was to assess the extent of the most widely used
382 formalisms in the literature for reproducing experimental pH-dependent distribution coefficients
383 in small molecules. To achieve this task, we built a database which consists of experimental pK_a ,
384 $\log P_N$, $\log P_I^{app}$, and $\log D_{7.4}$ values reported by Avdeef.²⁹ In addition, we employed experimental
385 entries of 86 molecules from the work of Tsantili-Kakoulidou *et al.*, containing $\log D_{pH}$ values at
386 various pH for each molecule as an individual entry.⁵⁵ Molecules with $\log D_{pH}$ values measured in
387 the presence of background salt concentrations above 0.15 mol/L were discarded because the study
388 of the effect of external ions on lipophilicity is beyond the scope of our study. Thus, we obtained
389 225 entries (118 individual molecules) with 113 acids, 100 bases, and 12 zwitterions. The limited
390 number of data entries in our database lies in the lack of publicly available data of experimental
391 measurements of pH-dependent lipophilicity profiles of molecules, especially the limited number
392 of apparent ionic partition coefficient measurements in the literature.

393 The $\log D_{pH}$ was calculated using Eq. 1 and Eq. 2 for each molecule at their respective pH
394 values. Figure 4 shows the overall performance of each model by comparing the modeled values
395 with their respective experimental $\log D_{pH}$ values. As expected, most of the molecules whose

396 $\log D_{\text{pH}}$ values were measured under different pH conditions to 7.4, present the largest deviation
397 using Eq. 1 (see Figure 4a, red marks), with highly underestimated predictions. As a consequence,
398 Eq. 1 poorly predicts $\log D_{\text{pH}}$ values at extreme pH values. On the other hand, the predicted values
399 using the Eq. 2 are significantly better (see Figure 4b), reducing the RMSE by 0.48 $\log D$ units,
400 which represents an improvement of 55% in accuracy.

401



402

403 **Figure 4.** Evaluation of the computed $\log D_{\text{pH}}$ of our database compared with the experimental
404 values with (a) Eq. 1 and (b) Eq. 2. Rhomboids represent $\log D_{\text{pH}}$ when the pH is different of 7.4.
405 Red dots and rhomboids highlight compounds with deviations greater than 1.5 $\log D$ units.
406 Statistical parameters were calculated using the ‘Metrics’ package in R (R^2 = squared Pearson’s
407 correlation coefficient, RMSE = root mean squared error, MAE = mean absolute error, and , MSE
408 = mean signed error).

409

410 Table 2 shows the reduction of RMSE in $\log D$ units of each molecule type using Eq.2
411 instead of Eq.1. It is observed that our dataset shows a significant improvement (ca. 54 %) in its
412 performance when its distribution coefficient is modeled with $\log D_{\text{Eq.2}}$ (see Figure S5). Basic
413 molecules showed the greatest improvement, amounting to 66 %, whereas the acid ones 44 %.
414 Zwitterions also showed a significant improvement (ca. 30 %), even though these molecules can
415 have multiple ionic partition coefficients (cationic partitions P^+ , anionic partitions P^- , and
416 zwitterionic partitions P^z), which are not considered in the model $\log D_{\text{Eq.2}}$. These partitions can be

417 added by considering both acidic and basic pK_a into the thermodynamic equilibria.⁴⁵ Despite this,
418 the implementation of just one of the two P_I^{app} show an improvement in the lipophilic modelling
419 of zwitterions since the pH conditions favored one of these possible ionic species over the others.
420

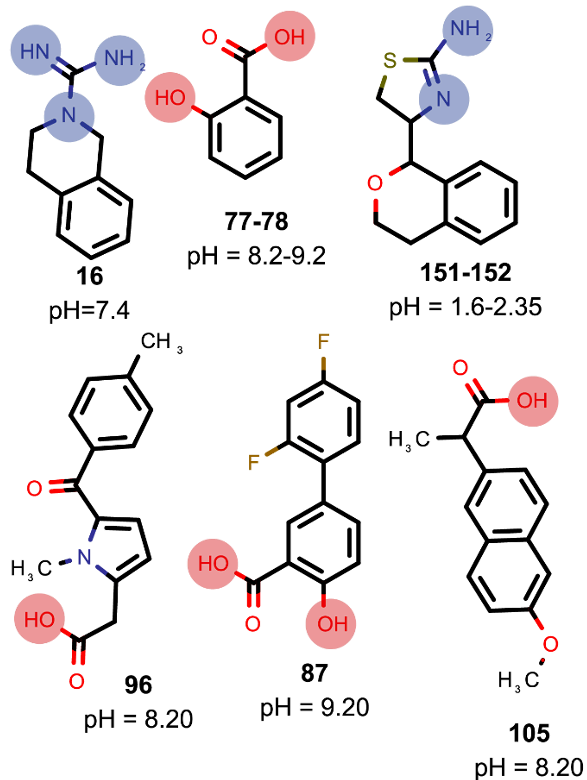
421 **Table 2.** Values of $\Delta RMSE$ for each type of molecule analyzed within our dataset by comparing
422 the modelled lipophilicities by $\log D_{Eq,1}$ and $\log D_{Eq,2}$ with their experimental values (Figure S1).

| Type | $\Delta RMSE^a$ |
|------------|-----------------|
| Acid | 0.30 |
| Base | 0.67 |
| Zwitterion | 0.38 |
| All | 0.48 |

423 ^a $\Delta RMSE = RMSE(\log D_{Eq,1}) - RMSE(\log D_{Eq,2})$

424

425 The molecules with the highest deviations in the prediction of experimental $\log D_{pH}$ using
426 $\log D_{Eq,1}$ are displayed in Figure 5. The chemical nature of the outliers is dominated by the presence
427 of ionic species because these compounds were experimentally measured under extreme pH
428 conditions. These deviations correspond to the theoretical frameworks of Eqs. 1 and 2. Thus, the
429 inclusion of the term P_I^{app} in Eq. 2 significantly corrects the prediction. Figure 5 also shows various
430 polyacids with pK_a values separated by more than four orders of magnitude, allowing us to analyze
431 the distribution coefficient using the most acidic pK_a . Bases with entries **16**, **151-152** have multiple
432 protonation sites, while acids with entries **77**, **78**, and **87** have two deprotonation sites. More
433 complex thermodynamic models can be considered for these molecules.⁴⁵ However, as mentioned
434 above, because of the separation of their pK_a , the consideration of P_I^{app} for the carboxylate species
435 with $\log D_{Eq,2}$ is enough to remarkably increase the accuracy of the lipophilicity prediction of these
436 compounds to extreme pH where one charged species can predominate over the others.



437

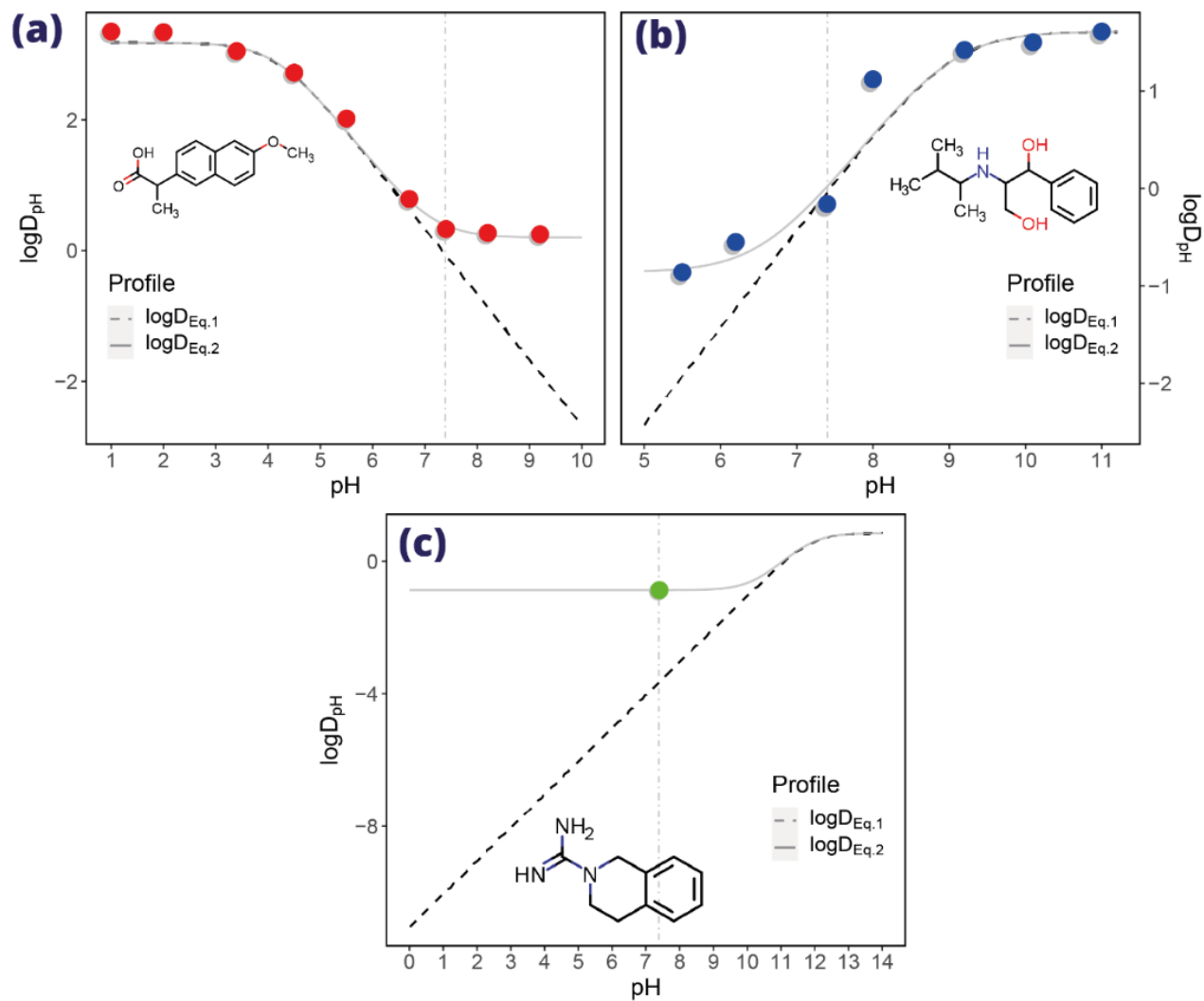
438 **Figure 5.** Representation of the molecules with the highest deviations in the prediction of the
 439 experimental $\log D_{\text{pH}}$ using $\log D_{\text{Eq1}}$. The protonation and deprotonation sites of each molecule are
 440 labeled in blue and red, respectively.

441

442 We also tested the performance of both formalisms by evaluating the entire pH-
 443 lipophilicity profile of individual molecules. To this end, an acidic (Naproxen) and two basic
 444 (compound '1774' from Ref 55 and Debrisoquine) examples were used. Naproxen and the
 445 compound '1774' (see Figure 6a-b) were selected because of the large amount of experimental
 446 data available in our database.⁵⁶ Therefore, it is better appreciated how the behavior of the
 447 experimental lipophilicity profiles fits more closely when evaluated using Eq. 2, particularly at
 448 extreme pH values. Additionally, at pH 7.4, the influence of apparent ionic partitioning can be
 449 observed depending on the chemical nature of the molecule. For instance, Debrisoquine in Figure
 450 6c represents a base with a very high pK_a . Hence, ionic species were more abundant at pH 7.4,
 451 which aligned more closely with the lipophilicity profile determined by Eq. 2. These results
 452 indicate that to reproduce the pH-dependent lipophilicity profiles of small molecules, it is

453 recommended to use Eq. 2, especially under pH conditions where ionic species are more
454 representative than neutral species.

455



456

457 **Figure 6.** Calculated pH-dependent lipophilic profiles of (a) acidic (Naproxen) and (b-c) basic
458 molecules (compound '1774' from Ref 55 and Debrisoquine) within our dataset. The dashed lines
459 represent the $\log D_{\text{pH}}$ values calculated using Eq. 1, and the solid line represents the values
460 calculated using Eq. 2. The dots represent experimental $\log D_{\text{pH}}$ values. A dot-dashed vertical line
461 was placed at pH 7.4.

462

463 **Use of pH-dependent distribution coefficients in medicinal, food, and environmental**
464 **chemistry.**

465
466 Although the distribution coefficient in solvent systems represents only a mimetic for many
467 biological and physicochemical processes, its relevance and successful application in several life
468 sciences fields is undeniable. In this regard, we further investigated the repercussions of the
469 apparent ionic partition of molecules in the applied parameters and metrics where lipophilicity is
470 relevant.

471 First, owing to the availability of experimental values for pH-dependent distribution
472 coefficients in the olive oil/water system for two bioantioxidants⁹⁶, we simulated the $\log D_{4.5}$ for
473 these phenolic acids using Eq. 1 and Eq. 2. Table 3 shows that Equation 2 fits best with gallic acid.
474 On the other hand, for caffeic acid, an appreciable error was observed using both formalisms,
475 amounting to almost 1 $\log D$ unit.

476
477 **Table 3.** Experimental and modeled distribution coefficients for the two bioantioxidants to
478 a pH of 4.5 in the olive oil/water system using Eq. 1 and Eq. 2.

| Compound | Experimental values (Ref 96) | | | | Calculated $\log D_{4.5}$ ($\Delta \log D^a$) | |
|--------------|------------------------------|--------|------------------|----------------|---|-------------|
| | $\log P_{oil/water}$ | pK_a | $\log P_I^{app}$ | $\log D_{4.5}$ | Eq. 1 | Eq. 2 |
| Gallic acid | 2.97 | 4.40 | 2.34 | 2.70 | 2.62 (-0.08) | 2.73 (0.03) |
| Caffeic acid | 3.26 | 4.54 | 1.70 | 2.04 | 2.98 (0.94) | 2.99 (0.95) |

479 ^a $\Delta \log D = (\text{calc} - \text{exp})$

480

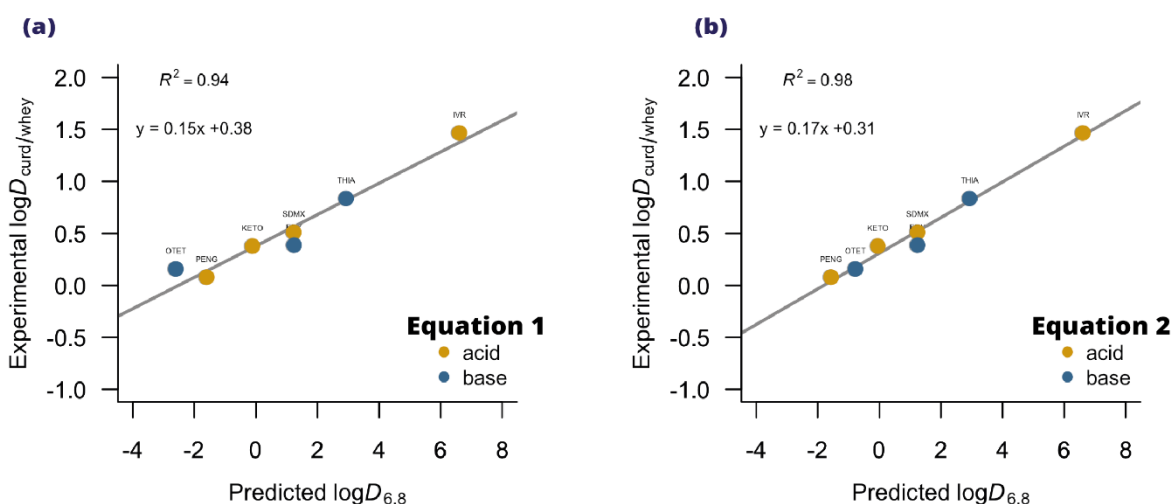
481

482

483

484 Second, previous studies have shown that the distribution of spiked drugs between milk
485 fractions, for example, the curd/whey system, to a pH of 6.8, can be properly mimicked through
486 the *n*-octanol/water distribution coefficient ($\log D_{6.8}$) using Eq. 1 (see Figure 7a).^{98,99} Despite the
487 excellent results obtained using Eq. 1, we were interested in investigating whether the use of Eq.
488 2 further improves the observed model (see Figure 7b).

489



490

491 **Figure 7.** Comparison between *n*-octanol/water $\log D_{6.8}$ using Eq. 1 (a) and Eq. 2 (b) and
492 the experimental distribution of spiked drugs between the curd/whey milk fractions. Drugs that
493 represent each acronym in the plot are listed in Table 4.

494

495

496

497

498

499

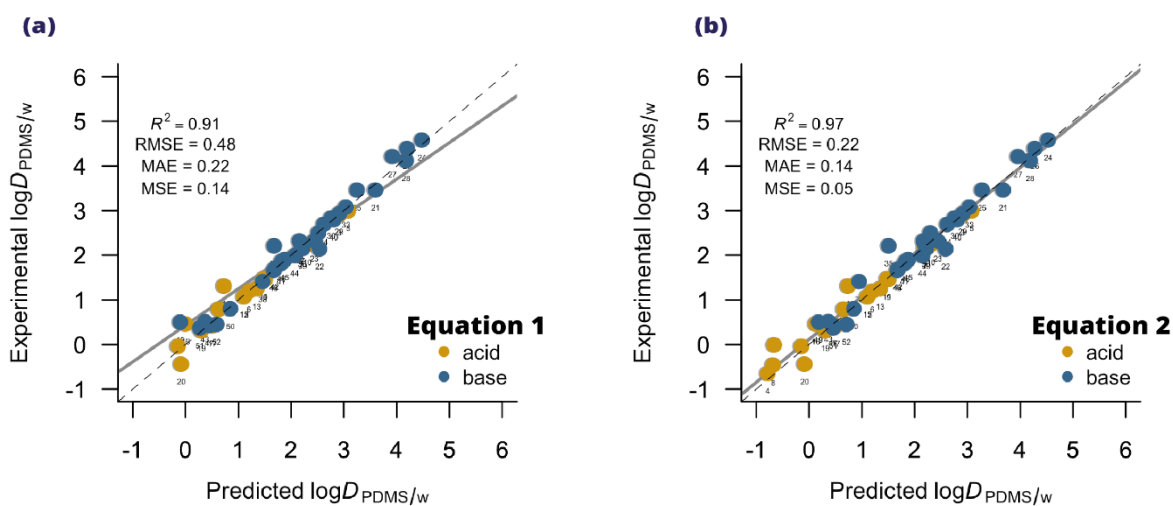
500 Table 4 reports the data used in the previous report^{98,99} as well as the $\log P_1^{\text{app}}$ of the tested
 501 molecules picked up from experimental measurements from the literature when available;
 502 otherwise, this parameter was simulated using predictive software such as ChemAxon.^{94,100} The
 503 predicted *n*-octanol/water distribution coefficients for spiked drugs to pH = 6.8 using Eq. 2 (see
 504 Fig 7, right) improved the correlation by 4% and showed an improved linear regression between
 505 both descriptors (see the linear equations in Fig. 7). In more detail, the improvement resides
 506 precisely in two compounds that had an experimental value of $\log P_1^{\text{app}}$ (ketoprofen and
 507 oxytetracycline). This observation highlights the importance of experimental measurements of
 508 ionic partitions but also calls for more experimental work focused on these issues, taking into
 509 account the scarce values available in the literature and the difficulty of finding them in public
 510 databases.

511 **Table 4.** Experimental distribution of spiked drugs between curd/whey milk fractions and
 512 predicted *n*-octanol/water distribution coefficients for spiked drugs at pH = 6.8 using Eq. 1 and
 513 Eq. 2.

| Compound | Original data from Refs. 98 and 99 | | | $\log P_1^{\text{app}}$ | Calculated $\log D_{6.8}$ | |
|----------------------------|------------------------------------|--------|-------------------------------------|-------------------------|---------------------------|-------|
| | $\log P_N$ | pK_a | $\log D$ curd/whey (pH = 6.8) | | Eq. 1 | Eq. 2 |
| Penicillin G (PENG) | 1.67 | 3.53 | 0.08 | -2.75 | -1.60 | -1.57 |
| Sulfadimethoxine (SDMX) | 1.48 | 6.91 | 0.51 | 0.12 | 1.23 | 1.24 |
| Ketoprofen (KETO) | 2.81 | 3.88 | 0.38 | -0.95 ^a | -0.11 | -0.05 |
| Ivermectin B1a (IVR) | 6.61 | 12.47 | 1.47 | 1.55 | 6.61 | 6.61 |
| Oxytetracycline (OTET) | -1.60 | 7.75 | 0.16 | -0.74 ^b | -2.60 | -0.78 |
| Erythromycin A (ERY) | 2.83 | 8.38 | 0.39 | -0.89 ^b | 1.24 | 1.24 |
| Thiabendazole (THIA) | 2.93 | 4.08 | 0.84 | 1.28 | 2.93 | 2.93 |

514 ^aExperimental data reported in our database in Ref 56. ^b Experimental data reported in Ref. 99.

515 Additionally, in environmental chemistry research, passive equilibrium sampling of
516 dissolved contaminants in water has been studied using polymer polydimethylsiloxane (PDMS) as
517 an absorbent phase. This hydrophobic passive sampler can extract ionizable compounds from
518 sediments and suspended particulate matter in a pH-dependent manner.¹⁰¹ The authors tested the
519 partitioning of ionizable compounds between PDMS and water ($\log D_{\text{PDMS/w}}$) at different pH
520 values. Thus, $\log D_{\text{PDMS/w}}$ measurements at extreme pH were considered as $\log P_{\text{N}}$ or $\log P_{\text{I}}^{\text{app}}$
521 depending on the acidic or basic nature of each molecule. These values were used to calculate
522 $\log D_{\text{Eq.1}}$ and $\log D_{\text{Eq.2}}$ in order to reproduce the experimental $\log D_{\text{PDMS/water}}$ measurements to a pH
523 of 7.4 (see Supporting Information, *environmental.csv*).



524
525 **Figure 8.** Comparison between PDMS and water $\log D_{7.4}$ using Eq. 1 (left), and Eq. 2 (right), and
526 the experimental $\log D_{\text{PDMS/w}}$ for a series of 52 ionizable compounds.

527
528
529
530
531

532 Figure 8 shows that the consideration of the apparent ionic partition significantly improved
533 the correlation between the experimental and predicted values, where the RMSE decreased in 0.26
534 $\log D$ units, representing an appreciable improvement of 54 %. This application demonstrated that
535 the thermodynamic equilibrium derived in Eq. 2 applies to partitions other than the n-octanol/water
536 system, such as in the PDMS/water phases. The presence of some free silanol groups, combined
537 with the highly hydrophobic polymeric chain might create a suitable environment for ionized-
538 organic species in the PDMS phase through a combination of hydrogen bonds from the terminal -
539 OH groups and dispersion non-covalent forces from the polymeric chains.^{101,132}

540 Finally, an important metric that has been increasingly applied in drug discovery and
541 medicinal chemistry lead optimization endeavors is the lipophilic efficiency (LipE, see Eq. 6).
542 LipE relates the binding affinity and lipophilicity of a compound, which creates a significant
543 parameter for estimating druglikeness.¹⁰² A proper interval of lipophilicity at physiological pH
544 ($\log D_{7.4}$), usually between 1 and 3, underpins the desired ADME properties and dose; therefore,
545 improving potency without excessively increasing lipophilicity is of vital importance in drug
546 discovery optimization programs. Table 5 compiles the LipE values obtained in the literature using
547 strictly experimental *pAct* lipophilicity at physiological pH, the latter due to the impossibility of
548 finding data at other pH values in the literature.

549

550

551

552

553

554

555

556

557

558

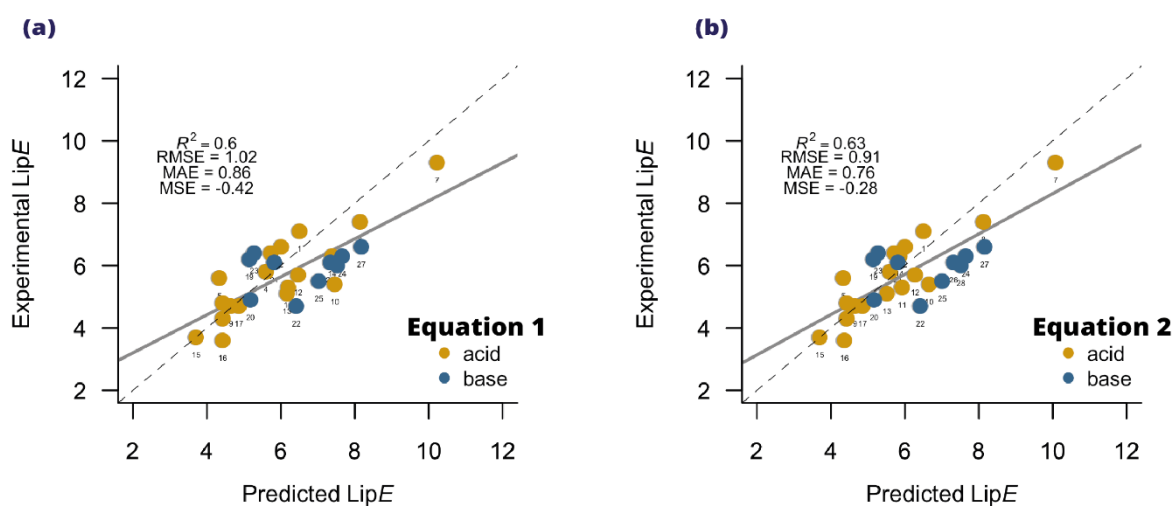
559

560 **Table 5.** Experimental lipophilic efficiency (LipE) for compounds reported in the literature and
 561 predicted lipophilic efficiency at pH = 7.4 using Eq. 1 and Eq. 2. The number representing the
 562 compounds in each original work is placed in column 'Id'.

| Compound | Id (Ref) | Exp. pAct | Exp. logD _{7.4} | Exp. LipE | logP _N | pK _a | logP ₁ ^{app} | LipE Eq. 1 | LipE Eq. 2 |
|----------|--------------------|-----------|--------------------------|-----------|-------------------|-----------------|----------------------------------|------------|------------|
| 1 | 6 (112) | 8.5 | 1.4 | 7.1 | 2.07 | 8.42 | -0.20 | 6.51 | 6.51 |
| 2 | 8 (112) | 7.9 | 1.3 | 6.6 | 1.88 | 8.42 | -0.39 | 6.01 | 6.01 |
| 3 | 9 (112) | 8.1 | 1.7 | 6.4 | 2.37 | 8.42 | 0.09 | 5.73 | 5.72 |
| 4 | 10 (112) | 7.8 | 2.1 | 5.8 | 2.29 | 8.42 | 0.03 | 5.60 | 5.60 |
| 5 | 11 (112) | 7.9 | 2.3 | 5.6 | 3.59 | 8.42 | 1.32 | 4.34 | 4.34 |
| 6 | 12 (112) | 7.8 | 3.5 | 4.3 | 3.43 | 8.42 | 1.16 | 4.43 | 4.43 |
| 7 | Rosuvastatin (109) | 9.0 | -0.3 | 9.3 | 1.90 | 4.27 | -1.63 | 10.23 | 10.08 |
| 8 | Pravastatin (109) | 7.1 | -0.2 | 7.4 | 2.18 | 4.20 | -2.41 | 8.15 | 8.13 |
| 9 | Fluvastatin (109) | 6.6 | 1.9 | 4.7 | 4.17 | 4.30 | 0.18 | 4.65 | 4.65 |
| 10 | 8 (108) | 8.4 | 3.0 | 5.4 | 4.50 | 3.84 | 1.67 | 7.46 | 6.66 |
| 11 | 9 (108) | 7.7 | 2.4 | 5.3 | 4.99 | 3.89 | 1.42 | 6.20 | 5.93 |
| 12 | 12 (108) | 8.2 | 2.5 | 5.7 | 5.01 | 4.07 | 1.41 | 6.47 | 6.28 |
| 13 | 13 (108) | 8.3 | 3.2 | 5.1 | 5.61 | 3.92 | 2.67 | 6.17 | 5.52 |
| 14 | 14 (108) | 8.3 | 2.0 | 6.3 | 4.23 | 4.08 | 2.43 | 7.39 | 5.86 |
| 15 | 8 (107) | 6.4 | 2.7 | 3.7 | 2.77 | 7.90 | 1.50 | 3.71 | 3.7 |
| 16 | 10 (107) | 5.9 | 2.3 | 3.6 | 2 | 7.00 | 0.75 | 4.43 | 4.37 |
| 17 | 11 (107) | 6.9 | 2.3 | 4.7 | 2.3 | 7.60 | 1.00 | 4.87 | 4.86 |
| 18 | 19 (107) | 4.8 | 0.1 | 4.8 | 0.43 | 9.50 | -2.00 | 4.43 | 4.43 |
| 19 | Indinavir (104) | 9.1 | 2.9 | 6.2 | 2.92 | 6.20 | -2.42 | 5.15 | 5.15 |
| 20 | 3 (104) | 8.0 | 3.2 | 4.9 | 4.5 | 8.97 | 0.77 | 5.18 | 5.17 |
| 21 | Crizotinib (104) | 8.1 | 1.9 | 6.1 | 4.28 | 9.40 | -0.38 | 5.82 | 5.81 |
| 22 | 8l (112) | 6.4 | 1.7 | 4.7 | 1.78 | 5.66 | -1.96 | 6.42 | 6.42 |
| 23 | 22 (102) | 7.0 | 3.6 | 6.4 | 4.24 | 8.87 | 0.43 | 5.28 | 5.28 |
| 24 | 7a (113) | 9.2 | 2.8 | 6.3 | 2.65 | 9.60 | -1.13 | 7.66 | 7.65 |
| 25 | 7j (113) | 5.6 | 1.0 | 5.5 | 3.38 | 9.70 | -0.41 | 7.03 | 7.01 |
| 26 | 7k (113) | 7.9 | 1.6 | 6.1 | 3.04 | 9.70 | -0.75 | 7.34 | 7.32 |
| 27 | 7m (113) | 6.8 | 2.1 | 6.6 | 2.22 | 9.70 | -1.57 | 8.18 | 8.16 |
| 28 | 7s (113) | 7.9 | 2.4 | 6.0 | 2.89 | 9.70 | -0.90 | 7.53 | 7.52 |

563

564 **Figure 9** shows the LipE simulated using Eq 1. and 2 compared to their experimental LipE
 565 values. Eq. 2 again shows favorable statistical parameters compared with Eq.1, improving the
 566 RSME in log D units by 11 %. Let us mention that the reproduction of LipE in both cases was not
 567 very satisfactory, this may be mainly due to the lack of experimental data of $\log P_N$ values for these
 568 compounds, but in particular to the use of simulated
 569 $\log P_1^{\text{app}}$ values. Although, tools such as ChemAxon have presented very good results in partition
 570 coefficient predictions, the molecules included here belong to novel compounds reported in
 571 medicinal chemistry articles with new chemical spaces that may impact the performance of
 572 predictive tools, especially in $\log P_1^{\text{app}}$ which to the best of our knowledge, this is the first work in
 573 reporting a free database for this parameter taken from several reports in the literature.



574
 575 **Figure 9.** Comparison of the lipophilic efficiency (LipE) using Eq. 1 (left), and Eq. 2 (right), and
 576 the experimental lipophilic efficiency for a series of 28 ionizable compounds.

577
 578
 579
 580
 581

582 To summarize, it can be noted that the inclusion of apparent ionic partitions can improve
583 metrics in the modeling of various metrics where lipophilicity is crucial to simulate biological or
584 artificial environments of higher complexity. Of special interest, we demonstrated that these
585 formalisms can be applied to systems beyond the classical *n*-octanol/water system. The
586 improvements studied ranged from 4 to 54 % in terms of RMSE ($\log D$ units) and depended mainly
587 on the pH at which the system was being simulated, pKa, and hydrophobicity of the molecules.

588

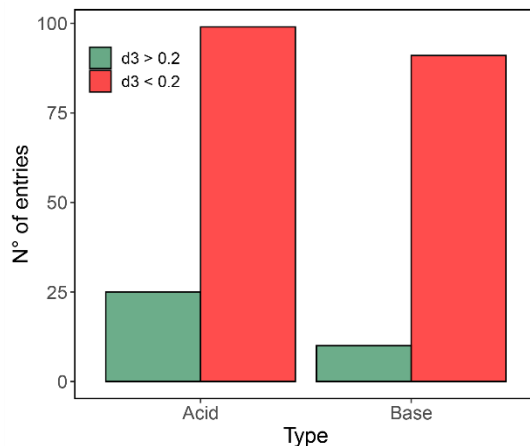
589 **Machine Learning models to guide the choice of pH-dependent lipophilicity profiles as a** 590 **function of molecular properties.**

591 The application of Eq. 1 offers significant advantages due to its simplicity of
592 implementation. However, the preceding sections emphasize the importance of employing Eq. 2
593 in various scenarios. It is important to note that the application of this formalism is constrained by
594 the availability of experimental data for partitioning of ionic species and simulations that ensure
595 adequate accuracy. This led us to propose a model that will help as a guide to discern cases in
596 which the simplified model represented in Eq. 1 can be used, or the consideration of the apparent
597 ionic partition is mandatory, as in Eq.2.

598 Consequently, one of the aims of this study was to develop a classification algorithm that
599 can differentiate whether the lipophilicity profile of a molecule can be better predicted with
600 $\log D_{\text{Eq.1}}$ or $\log D_{\text{Eq.2}}$. However, a significant number of entries indicate that both formalisms
601 compute a similar result compared to their experimental values by yielding d_3 values close to 0
602 (see Figure S6a). Let us note that we focus on the specific cases with a significant improvement
603 when the P_1^{app} of molecules is considered. Therefore, we decided to delimit the conditional d_3 ,
604 indicating that if a molecule exceeds a certain value of d_3 , it is important to consider its apparent
605 ionic partitioning to predict its lipophilicity. We tested d_3 values between 0.1-1 and picked the
606 optimal value based on two main parameters. First, considering that our set was small because we
607 used strictly experimental values in our database, we seek that the population of molecules that
608 best fit with $\log D_{\text{Eq.2}}$ should be at least 10 %. Then, there should be a sufficient number of
609 descriptors that have statistically proven divergence by WTT ($p < 0.05$). Thus, Machine Learning
610 algorithms have a larger number of parameters to create predictive models with higher accuracy.
611 Consequently, the delimiter ‘0.2’ showed an adequate balance between these two parameters and

612 was selected as our cut-off value (see Figure S6b). Molecules with values of $d_3 > 0.2$ showed an
613 improvement in lipophilicity modeling using Eq. 2. On the other hand, entries that had negative d_3
614 values or that fell into the range $0.2 < d_3 < 0$ were classified as molecules where the difference
615 between both models was negligible, and thus were classified as better fitted using $\log D_{\text{Eq.1}}$ due to
616 its easy implementation (it does not depend on $P_{\text{I}}^{\text{app}}$, resulting in less computational effort and
617 fewer experimental parameters). Higher thresholds significantly decreased the population in
618 $\log D_{\text{Eq.2}}$, while lower values reduced the structural divergence between molecules in $\log D_{\text{Eq.1}}$ and
619 $\log D_{\text{Eq.2}}$, making it more difficult to find descriptors that can differentiate between both
620 populations. The value '0.5' was also tested because a local maximum of descriptors with $p < 0.05$
621 was observed at this point (see Figure S6b). Furthermore, this value is of experimental interest,
622 because $\log P_{\text{N}}$ measurements of substances with different techniques tend to vary by less than 0.5
623 $\log P$ units (using the Shake-Flask method as a reference), which is considered as a parameter
624 indicating that the experimental techniques are not equivalent.³⁰ However, this extreme value and
625 the descriptors selected (see Table S14) showed poor performance in the ML models tested,
626 especially with External Set 1 (see Figure S7). This phenomenon can be explained because this d_3
627 delimiter has a very small $\log D_{\text{Eq.2}}$ population, thus the datasets are extremely unbalanced and the
628 robustness of the models is reduced. On the other hand, the accuracy of experimental methods,
629 even using different techniques, rounds at values less than 0.2 $\log P$ units.³⁰ Therefore, we
630 continued to train the ML models using the $d_3 > 0.2$ cut-off value to determine tendencies among
631 the selected descriptors via the feature selection methods and to evaluate the performance of the
632 ML algorithms.

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



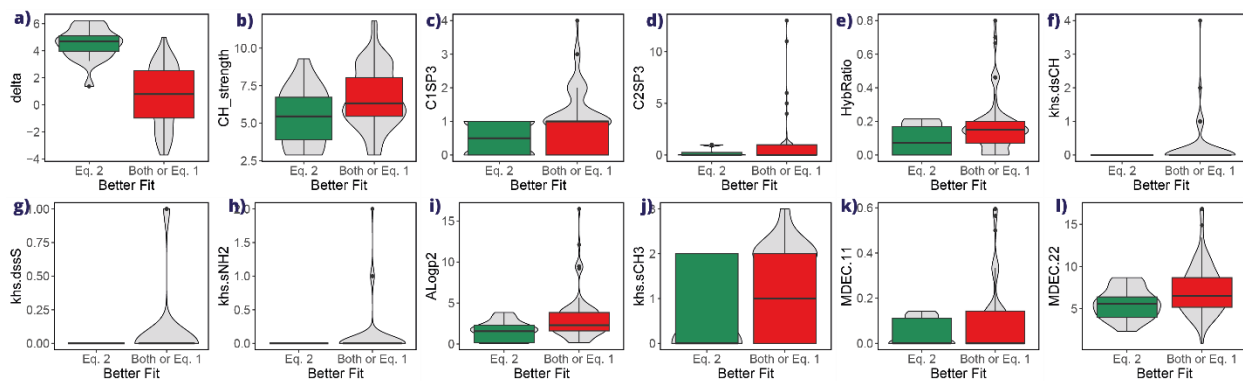
637
638 **Figure 10.** Distribution of acid and basic entries from our dataset as a function of their d_3 values.

639
640 We obtained several structural and physicochemical descriptors of the molecules to
641 identify considerable divergence between populations. First, our database was split into acids and
642 bases, and then into the training and test sets. The ‘*rcdk*’ package in R was used to look through
643 the descriptors, along with the *Jazzy* calculations of energies of hydration and hydrogen-bond
644 strengths and the experimental descriptors. The selected feature selection methods showed a wide
645 range of diverse descriptors (see Tables 1 and S1). We performed a Welch’s *t*-test on our
646 descriptors (WTT), which analyzes the divergence between populations relative to the variances
647 of the two groups.¹¹⁷ This test was selected over a Student’s *t*-test because of the divergence of
648 sample sizes (Figure 10) and variances between groups (Figure 11-12).¹¹⁸ The WTT descriptors
649 provided acceptable accuracies (see Figure S8).

650 An iterative feature selection method was also tested using the RFE model. The algorithm
651 achieved better performance when the 14 most important variables for acids and the nine most
652 important variables for bases were maintained. The importance of each descriptor posed by the
653 RFE is shown in Figure S3. Good results were obtained when these descriptors were implemented
654 during the training of the Machine Learning models. However, the accuracy decreased
655 significantly when the Test and External Set 1 were evaluated (see Figure S8c-d), indicating that
656 these descriptors did not generate a sufficiently robust model, or that the large number of chosen
657 descriptors (see Table S1) may overfit the data. Therefore, we selected WTT descriptors to analyze

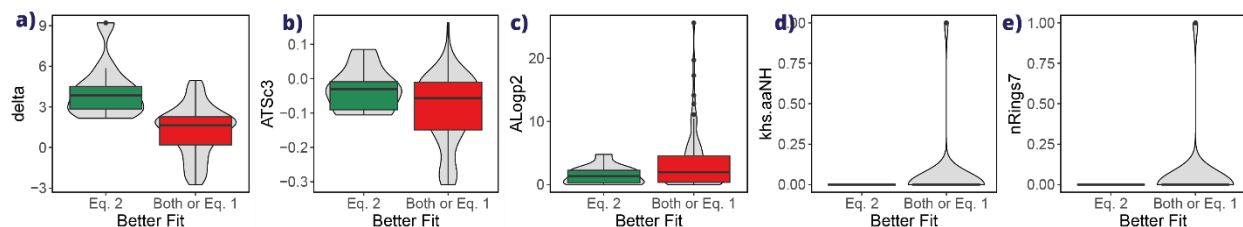
658 the tendencies of the molecules in each population and to evaluate the overall performance of the
659 Machine Learning algorithms that we developed.

660



661
662 **Figure 11.** Violin plots of the distribution of the acidic molecules in our dataset along the selected
663 descriptors for the acids ((a) *delta*, (b) *CH_strength*, (c) *C1SP3*, (d) *C2SP3*, (e) *HybRatio*, (f)
664 *khs.dsCH*, (g) *khs.dssS*, (h) *khs.sNH2*, (i) *Alogp2*, (j) *khs.sCH3*, (k) *MDEC.11*, and (l) *MDEC.22*).
665 Distributions are separated between acids and bases and classified by the binary operator $d_3 > 0.2$
666 (green) and $d_3 < 0.2$ (red).

667



668
669
670 **Figure 12.** Violin plots of the distribution of the acidic molecules in our dataset along the selected
671 descriptors for the bases (a) *delta*, (b) *ATSc3*, (c) *Alogp2*, (d) *khs.aanNH*, and (e) *nRings7*).
672 Distributions are separated between acids and bases and classified by the binary operator $d_3 > 0.2$
673 (green) and $d_3 < 0.2$ (red).

674

675

676 Figure 11 and 12 show the selected descriptors for acids and bases, respectively, used to
677 train our classification ML models. These descriptors showed statistically significant divergence
678 between the means of both populations among the 180 descriptors tested for acids and bases. Both
679 acidic and basic compounds showed significant differences in their means ($p < 0.05$ in WTT test)
680 for the *delta* and *Alogp2* descriptors (Table 1). The descriptor *delta* was calculated at the respective
681 pH of each entry for acids and bases. As expected, this descriptor proved to be the most important
682 in every test carried out in this regard (see Figures S3-S4), as it correlates with the prominence of
683 ionic species in both phases. Therefore, the apparent ionic partition becomes more significant for
684 entries with higher *delta* values (Figures 11a and 12a). This result is very promising, because
685 despite being an experimental descriptor, there are computational methods to determine pK_a that
686 include first-principles models^{133–136} as well as machine learning tools^{137,138}. Thus, the descriptor
687 *delta* can be automated and easily used to classify molecules according to the lipophilicity
688 formalisms analyzed here. In fact, the root-mean-square error (RMSE) between predicted pK_a
689 values using the software ChemAxon and experimental data in our database is just 0.58 log units
690 and the squared coefficient of determination (R^2) of 0.95 (see Figure S9)

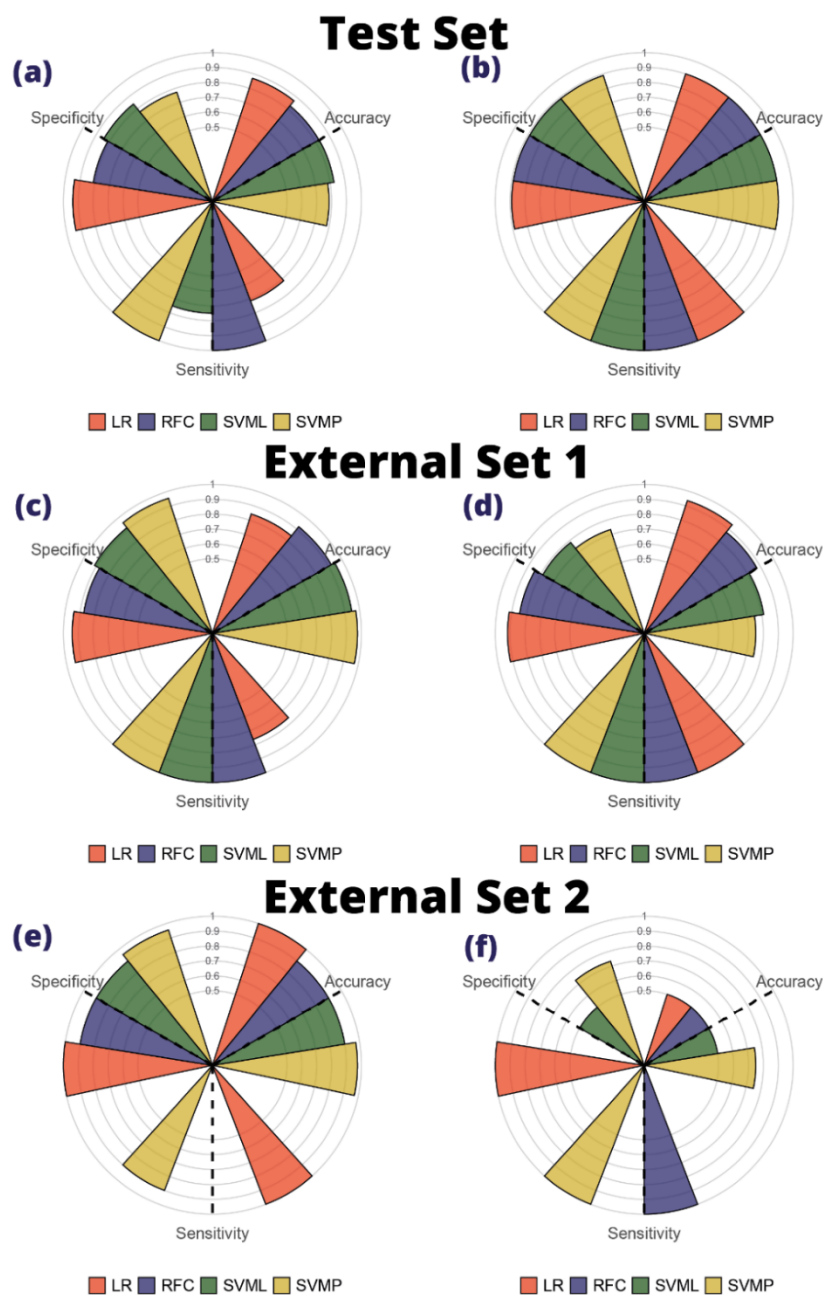
691 The *Alogp2* descriptor consists of a 3D-QSAR model by Ghose & Crippen (1986), which
692 predicts a square value of the $\log P_N$ value by analyzing the presence of 110 structural fragments
693 within the molecules.⁹⁵ Figure 11i and 12c show that molecules with hydrophobicity close to $\log P$
694 = 0 (with lower *Alogp2* values) tend to fit best with $\log D_{Eq,2}$. Although water and *n*-octanol are not
695 miscible, a small amount of water can dissolve in octanol at room temperature (~ 2.9 mol/kg).¹³⁹
696 These hydrophilic molecules might be dragged by the dissolved water to the octanol phase along
697 with the ionic species; thus, the apparent ionic partition would have a higher importance in these
698 molecules.

699 This affinity for water, at least for acidic compounds, was further demonstrated using the
700 *CH_strength* descriptor (Figure 11b). This descriptor, calculated by Jazzy, predicts the hydrogen-
701 bond donor strength in carbon atoms.¹¹⁶ The smaller *CH_strength* values indicate that for entries
702 with $d_3 < 0.2$, H-bond donors are not primarily found on carbons. Instead, they are found on other
703 more electronegative heteroatoms. Thus, by weakening the X-H covalent bonds through H-bonds,
704 the possibility of ionization of these species in both water and *n*-octanol increases. Figure 11e,k-l
705 presents other important descriptors for acidic compounds such as *MDEC.11*, *MDEC.22*, and
706 *HybRatio*. The *MDEC.11* and *MDEC.22* descriptor consists of a relationship between the number

707 of primary (*MDEC.11*) and secondary (*MDEC.22*) carbons in the molecule and the squared
708 average atomic distance between these atoms.¹²¹ Whereas, *HybRatio* is the number of sp³-C atoms
709 compared to the sum of sp³ and sp² C atoms. Eq. 2 works better for acidic substances with low
710 values of these descriptors, which considers together the values of *Alogp2*, allows us to infer that
711 small and rigid ionizable molecules with instaurations or aromatic systems need considering the
712 P_1^{app} to obtain an accurate prediction of $\log D_{\text{pH}}$.

713 Similarly, for basic compounds, higher values of the *ATSc3* descriptor are associated with
714 the consideration of the P_1^{app} for modeling pH-dependent lipophilic profiles on basic molecules
715 (Figure 12b). This descriptor is related to the high molecular polarizability, which agrees with the
716 pattern of small molecules in the presence of polar atoms such as nitrogen. Therefore, the apparent
717 ionic partition effect should be considered for these small, rigid, and unsaturated molecules, which
718 present a significant proportion of ionic species in the aqueous phase. It has been previously shown
719 that the P_1^{app} of molecules may mechanistically occur via a simple ion-transfer reaction.¹⁴⁰ Thus,
720 it is more plausible that small and compact molecules have a more prominent P_1^{app} because of the
721 lower energetic cost of transferring to the cavity of the ion they replace.

722 After establishing a distinct division between the two populations and applying an
723 appropriate feature selection method, Models A and B (see Figure 3) were trained using the logistic
724 regression (LR), random forest classification (RFC), and support vector machine (SVM)
725 algorithms. A training set for acidic and basic molecules was used for each model and was
726 evaluated using a test set consisting of 20% of the population. In addition, the two external sets
727 were validated using the experimental data obtained by Disdier et al. (External Set 1)⁴⁵ and
728 Fauchère and Pliška (External Set 2).¹²⁶ Predictions were made to determine which formalism best
729 modeled the lipophilicity of the inputs, and the results were collected in confusion matrices. The
730 performance of each marker was evaluated by calculating its accuracy, specificity, and sensitivity.
731 Figure 13 shows the results of the calculations of the four algorithms for the test and external sets
732 of acidic and basic molecules.



733

734

735 **Figure 13.** Accuracy, sensitivity, and specificity of each ML model evaluated in this study for
 736 acidic (a,c,e) and basic (b,d,f) entries within the test and external sets by defining our populations
 737 with the conditional $d_3 > 0.2$. Descriptors were selected using the WTT method. Accuracies,
 738 sensitivities, and specificities were calculated with Eqs. 10-12 based on the results of each
 739 confusion matrix (Tables S2-S13)

740 It is observed that most of the calculated accuracies for our test set have high values
741 (between 0.8 and 0.95), denoting that these classification models manage to distinguish relatively
742 well which molecules best fit with $\log D_{\text{Eq.1}}$ and $\log D_{\text{Eq.2}}$. However, the sensitivity decreased in the
743 test set of acidic molecules, indicating that the models had difficulties in detecting molecules that
744 fit $\log D_{\text{Eq.2}}$ (Figure 13a). External Set 1 exhibited good performance, with all models showing
745 similar accuracies, sensitivities, and specificities to those evaluated in the Test Set (Figure 13c-d).
746 Additionally, External Set 1 mainly comprises more hydrophobic molecules than our dataset, as
747 most molecules have $\log D_{\text{pH}}$ values < 0 (Figure S10). This demonstrates that our models exhibit
748 high robustness, even when dealing with species belonging to slightly different chemical spaces.
749 External Set 2, associated with capped amino acids as reported by Fauchère and Pliška¹²⁶, obtained
750 divergent results. On the one hand, the pH-dependent values of N-Acetyl-L-tyrosine amide were
751 predicted with excellent metrics, especially using the LR and SVMP models, because our training
752 set had a representative number of molecules with phenolic groups. On the other hand, in the case
753 of N-acetyl-L-histidine amide, the results were very poor due, at least in part, to the fact that our
754 set had few bases in relation to the acids that best fit Eq. 2, mainly because there was no imidazole
755 fragments present in our set of bases, thus limiting the performance of our models.

756
757
758
759
760
761
762
763
764
765
766
767

768 **Conclusions**

769 Lipophilicity is undoubtedly the most widely used and important descriptor in the early
770 stages of drug discovery and development. Additionally, it is a crucial descriptor in substance risk
771 assessment and in areas such as adsorption in materials, catalysts, food chemistry, and
772 computational biology. There are multiple tools to determine this descriptor, mainly for neutral
773 molecules ($\log P_N$). For substances with ionizable groups, two formalisms are commonly used to
774 determine the distribution coefficient ($\log D_{\text{pH}}$), being the simplest pH correction model is the most
775 widely used. However, previous studies carried out on specific and small molecule sets
776 recommend considering the effect of the apparent ionic compounds (P_I^{app}) because it has a negative
777 impact on the accuracy of computing lipophilic profiles when charged species or related species
778 are ignored. Our study, which was based on a larger amount of data and strictly on experimental
779 values, validated the observations presented in previous studies. We have also evidenced the
780 impact of P_I^{app} on the prediction of both the experimental lipophilicity profiles of small molecules
781 and experimental lipophilicity-based applications and metrics such as lipophilic efficiency (LipE),
782 distribution of spiked drugs in milk products, and pH-dependent partition of water contaminants
783 in synthetic passive samples such as silicones. Our findings show that better predictions are
784 obtained by considering the apparent ionic partition, whereas ignoring its contribution can lead to
785 inadequate experimental simplifications and/or computational predictions.

786 Finally, we developed machine learning algorithms using logistic regression, random forest
787 classification, and support vector machine models to determine from molecular structures in which
788 cases the P_I^{app} should be considered. The results indicate that small, rigid, and unsaturated
789 molecules with $\log P_N$ close to zero, which represent a significant proportion of ionic species in the
790 aqueous phase, are better modeled using the formalism that takes into account the apparent ionic
791 compounds (P_I^{app}).

792 Although we are aware of the molecular complexity of the species that can be included in
793 the computational determination of the apparent ionic partition (P_I^{app}), parameterization or training
794 of models using experimental values of P_I^{app} can help alleviate the restricted application of
795 formalisms that include this effect. Finally, our findings can serve as guidance to the scientific
796 community working in early-stage drug design, food, and environmental chemistry who deal with

797 ionizable molecules, to determine a priori which lipophilicity profile should be used depending on
798 the structure of a substance in research efforts. Future studies will address the influence played by
799 the apparent ionic partition (P_1^{app}) on the pH-dependent lipophilic profiles in more complex
800 systems such as zwitterionic and peptides.

801 **Data availability**

802 All codes are hosted at: https://github.com/cbio3lab/Lip_profiles. The database constructed in
803 this work can be consulted in the reference 56.

804 **Author contributions**

805 William J. Zamora: conceptualization, methodology, validation, data curation, writing – original
806 draft, writing – review & editing. Esteban Bertsch: methodology, validation, writing – review &
807 editing. Sebastián Suñer: methodology, validation, writing – review & editing. Silvana Pinheiro:
808 conceptualization, writing – review & editing

809 **Conflicts of interest**

810 There are no conflicts to declare.

811

812 **Acknowledgments**

813 The authors thank the Vice Chancellor for Research of the University of Costa Rica for its support
814 work via the research projects 115-C2-126 and 908-C3-610. We also thank Dr. Antonio Viayna of
815 the University of Barcelona for his comments for the improvement of this manuscript.

816

817

818

819

820 References

- 821 (1) Waring, M. J. Lipophilicity in Drug Discovery. *Expert Opinion on Drug Discovery* **2010**, *5*
822 (3), 235–248. <https://doi.org/10.1517/17460441003605098>.
- 823 (2) Lewis, D. F. V.; Jacobs, M. N.; Dickins, M. Compound Lipophilicity for Substrate Binding to
824 Human P450s in Drug Metabolism. *Drug Discovery Today* **2004**, *9* (12), 530–537.
825 [https://doi.org/10.1016/S1359-6446\(04\)03115-0](https://doi.org/10.1016/S1359-6446(04)03115-0).
- 826 (3) Chatzopoulou, M.; Emer, E.; Lecci, C.; Rowley, J. A.; Casagrande, A. S.; Moir, L.; Squire, S.
827 E.; Davies, S. G.; Harriman, S.; Wynne, G. M.; Wilson, F. X.; Davies, K. E.; Russell, A. J.
828 Decreasing HepG2 Cytotoxicity by Lowering the Lipophilicity of
829 Benzo[d]Oxazolephosphinate Ester Utrophin Modulators. *ACS Medicinal Chemistry Letters*
830 **2020**, *11* (12), 2421–2427.
831 <https://doi.org/10.1021/ACSMEDCHEMLETT.0C00405/ASSET/IMAGES/LARGE/ML0C0>
832 [0405_0003.JPEG](https://doi.org/10.1021/ACSMEDCHEMLETT.0C00405/ASSET/IMAGES/LARGE/ML0C0405_0003.JPEG).
- 833 (4) Miller, M. M.; Wasik, S. P.; Huang, G. L.; Shlu, W. Y.; Mackay, D. Relationships between
834 Octanol-Water Partition Coefficient and Aqueous Solubility. *Environmental Science and*
835 *Technology* **1985**, *19* (6), 522–529.
836 https://doi.org/10.1021/ES00136A007/ASSET/ES00136A007.FP.PNG_V03.
- 837 (5) Soliman, K.; Grimm, F.; Wurm, C. A.; Egner, A. Predicting the Membrane Permeability of
838 Organic Fluorescent Probes by the Deep Neural Network Based Lipophilicity Descriptor
839 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)
840 [021-86460-3](https://doi.org/10.1038/S41598-021-86460-3).
- 841 (6) Esser, H. O. A Review of the Correlation between Physicochemical Properties and
842 Bioaccumulation. *Pesticide Science* **1986**, *17* (3), 265–276.
843 <https://doi.org/10.1002/PS.2780170310>.
- 844 (7) Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and Computational
845 Approaches to Estimate Solubility and Permeability in Drug Discovery and Development
846 Settings. *Advanced Drug Delivery Reviews* **2001**, *46* (1–3), 3–26.
847 [https://doi.org/10.1016/S0169-409X\(00\)00129-0](https://doi.org/10.1016/S0169-409X(00)00129-0).
- 848 (8) Veber, D. F.; Johnson, S. R.; Cheng, H. Y.; Smith, B. R.; Ward, K. W.; Kopple, K. D.
849 Molecular Properties That Influence the Oral Bioavailability of Drug Candidates. *Journal of*
850 *Medicinal Chemistry* **2002**, *45* (12), 2615–2623.
851 https://doi.org/10.1021/JM020017N/SUPPL_FILE/JM020017N_S.PDF.
- 852 (9) Leo, A.; Hansch, C.; Church, C. Comparison of Parameters Currently Used in the Study of
853 Structure-Activity Relationships. *Journal of Medicinal Chemistry* **1969**, *12* (5), 766–771.
854 https://doi.org/10.1021/JM00305A010/ASSET/JM00305A010.FP.PNG_V03.
- 855 (10) Leo, A.; Hansch, C.; Elkins, D. Partition Coefficients and Their Uses. *Chemical Reviews*
856 **1971**, *71* (6), 525–616.
857 https://doi.org/10.1021/CR60274A001/ASSET/CR60274A001.FP.PNG_V03.
- 858 (11) Kerns, E. H. High Throughput Physicochemical Profiling for Drug Discovery. *Journal of*
859 *Pharmaceutical Sciences* **2001**, *90* (11), 1838–1858. <https://doi.org/10.1002/JPS.1134>.

- 860 (12) Liu, X.; Tu, M.; Kelly, R. S.; Chen, C.; Smith, B. J. Development of a Computational
861 Approach to Predict Blood-Brain Barrier Permeability. *Drug metabolism and disposition: the*
862 *biological fate of chemicals* **2004**, 32 (1), 132–139. <https://doi.org/10.1124/DMD.32.1.132>.
- 863 (13) Colmenarejo, G. In Silico Prediction of Drug-Binding Strengths to Human Serum
864 Albumin. *Medicinal Research Reviews* **2003**, 23 (3), 275–301.
865 <https://doi.org/10.1002/MED.10039>.
- 866 (14) Zamora, W. J.; Curutchet, C.; Campanera, J. M.; Luque, F. J. Prediction of PH-Dependent
867 Hydrophobic Profiles of Small Molecules from Miertus-Scrocco-Tomasi Continuum
868 Solvation Calculations. *Journal of Physical Chemistry B* **2017**, 121 (42), 9868–9880.
869 https://doi.org/10.1021/ACS.JPCB.7B08311/SUPPL_FILE/JP7B08311_LIVESLIDES.MP4.
- 870 (15) Iglesias, V.; Pintado-Grima, C.; Santos, J.; Fornt, M.; Ventura, S. Prediction of the Effect
871 of PH on the Aggregation and Conditional Folding of Intrinsically
872 Disordered Proteins. *In Data Mining Techniques for the Life Sciences*; Carugo, O., Eisenhaber,
873 F., Eds.; Methods in Molecular Biology; Springer US: New York, NY, 2022; pp 197–211.
874 https://doi.org/10.1007/978-1-0716-2095-3_8.
- 876 (16) Oeller, M.; Kang, R.; Bell, R.; Ausserwöger, H.; Sormanni, P.; Vendruscolo, M. Sequence-
877 Based Prediction of PH-Dependent Protein Solubility Using CamSol. *Briefings in*
878 *Bioinformatics* **2023**, 24 (2), bbad004. <https://doi.org/10.1093/bib/bbad004>.
- 879 (17) Porto, W. F.; Ferreira, K. C. V.; Ribeiro, S. M.; Franco, O. L. Sense the Moment: A Highly
880 Sensitive Antimicrobial Activity Predictor Based on Hydrophobic Moment. *Biochim Biophys*
881 *Acta Gen Subj* **2022**, 1866 (3), 130070. <https://doi.org/10.1016/j.bbagen.2021.130070>.
- 882 (18) Zamora, W. J.; De Souza, S.; Separovic, F.; Luque, Fco. J. Insights into the Effect of the
883 Membrane Environment on the Three-Dimensional Structure-Function Relationship of
884 Antimicrobial Peptides. *Biophysical Journal* **2020**, 118 (3, Supplement 1), 236a.
885 <https://doi.org/10.1016/j.bpj.2019.11.1394>.
- 886 (19) Simm, S.; Einloft, J.; Mirus, O.; Schleiff, E. 50 Years of Amino Acid Hydrophobicity
887 Scales: Revisiting the Capacity for Peptide Classification. *Biological Research* **2016**, 49 (1),
888 1–19. <https://doi.org/10.1186/S40659-016-0092-5/FIGURES/8>.
- 889 (20) Zamora, W. J.; Campanera, J. M.; Luque, F. J. Development of a Structure-Based, PH-
890 Dependent Lipophilicity Scale of Amino Acids from Continuum Solvation Calculations.
891 *Journal of Physical Chemistry Letters* **2019**, 10 (4), 883–889.
892 https://doi.org/10.1021/ACS.JPCLETT.9B00028/SUPPL_FILE/JZ9B00028_LIVESLIDES.MP4.
- 894 (21) Ingram, T.; Richter, U.; Mehling, T.; Smirnova, I. Modelling of PH Dependent N-
895 Octanol/Water Partition Coefficients of Ionizable Pharmaceuticals. *Fluid Phase Equilibria*
896 **2011**, 305 (2), 197–203. <https://doi.org/10.1016/j.fluid.2011.04.006>.
- 897 (22) Chen, C.-S.; Lin, S.-T. Prediction of PH Effect on the Octanol–Water Partition Coefficient
898 of Ionizable Pharmaceuticals. *Ind. Eng. Chem. Res.* **2016**, 55 (34), 9284–9294.
899 <https://doi.org/10.1021/acs.iecr.6b02040>.
- 900 (23) Westall, J. C.; Leuenberger, C.; Schwarzenbach, R. P. Influence of PH and Ionic Strength
901 on the Aqueous-Nonaqueous Distribution of Chlorinated Phenols. *Environmental Science and*

- 902 *Technology* **1985**, *19* (2), 193–198.
903 https://doi.org/10.1021/ES00132A014/ASSET/ES00132A014.FP.PNG_V03.
- 904 (24) Xing, L.; Glen, R. C. Novel Methods for the Prediction of LogP, Pka, and LogD. *Journal*
905 *of Chemical Information and Computer Sciences* **2002**, *42* (4), 796–805.
906 <https://doi.org/10.1021/CI010315D/ASSET/IMAGES/LARGE/CI010315DF00006.JPEG>.
- 907 (25) Tetko, I. V.; Poda, G. I. Application of ALOGPS 2.1 to Predict Log D Distribution
908 Coefficient for Pfizer Proprietary Compounds. *Journal of Medicinal Chemistry* **2004**, *47* (23),
909 5601–5604.
910 <https://doi.org/10.1021/JM049509L/ASSET/IMAGES/LARGE/JM049509LF00001.JPEG>.
- 911 (26) Livingstone, D. Theoretical Property Predictions. *Current Topics in Medicinal Chemistry*
912 **2005**, *3* (10), 1171–1192. <https://doi.org/10.2174/1568026033452078>.
- 913 (27) Bannan, C. C.; Calabró, G.; Kyu, D. Y.; Mobley, D. L. Calculating Partition Coefficients
914 of Small Molecules in Octanol/Water and Cyclohexane/Water. *J. Chem. Theory Comput.*
915 **2016**, *12* (8), 4015–4024. <https://doi.org/10.1021/acs.jctc.6b00449>.
- 916 (28) Bergazin, T. D.; Tielker, N.; Zhang, Y.; Mao, J.; Gunner, M. R.; Francisco, K.; Ballatore,
917 C.; Kast, S. M.; Mobley, D. L. Evaluation of Log P, PK a, and Log D Predictions from the
918 SAMPL7 Blind Challenge. *Journal of Computer-Aided Molecular Design* **2021**, *35* (7), 771–
919 802. <https://doi.org/10.1007/s10822-021-00397-3>.
- 920 (29) Avdeef, A. *Absorption and Drug Development.*; John Wiley & Sons, 2012.
- 921 (30) Port, A.; Bordas, M.; Enrech, R.; Pascual, R.; Rosés, M.; Ràfols, C.; Subirats, X.; Bosch,
922 E. Critical Comparison of Shake-Flask, Potentiometric and Chromatographic Methods for
923 Lipophilicity Evaluation (Log Po/w) of Neutral, Acidic, Basic, Amphoteric, and Zwitterionic
924 Drugs. *European Journal of Pharmaceutical Sciences* **2018**, *122*, 331–340.
925 <https://doi.org/10.1016/J.EJPS.2018.07.010>.
- 926 (31) Austin, R. P.; Barton, P.; Davis, A. M.; Manners, C. N.; Stansfield, M. C. The Effect of
927 Ionic Strength on Liposome–Buffer and 1-Octanol–Buffer Distribution Coefficients. *Journal*
928 *of Pharmaceutical Sciences* **1998**, *87* (5), 599–607. <https://doi.org/10.1021/js9703481>.
- 929 (32) Jain, P.; Kumar, A. Concentration-Dependent Apparent Partition Coefficients of Ionic
930 Liquids Possessing Ethyl- and Bi-Sulphate Anions. *Phys. Chem. Chem. Phys.* **2015**, *18* (2),
931 1105–1113. <https://doi.org/10.1039/C5CP06611E>.
- 932 (33) Jafvert, C. T.; Westall, J. C.; Grieder, E.; Schwarzenbach, R. P. Distribution of
933 Hydrophobic Ionogenic Organic Compounds between Octanol and Water: Organic Acids.
934 *Environ. Sci. Technol.* **1990**, *24* (12), 1795–1803. <https://doi.org/10.1021/es00082a002>.
- 935 (34) Austin, R. P.; Davis, A. M.; Manners, C. N. Partitioning of Ionizing Molecules between
936 Aqueous Buffers and Phospholipid Vesicles. *Journal of Pharmaceutical Sciences* **1995**, *84*
937 (10), 1180–1183. <https://doi.org/10.1002/jps.2600841008>.
- 938 (35) Takács-Novák, K.; Szász, G. Ion-Pair Partition of Quaternary Ammonium Drugs: The
939 Influence of Counter Ions of Different Lipophilicity, Size, and Flexibility. *Pharm Res* **1999**,
940 *16* (10), 1633–1638. <https://doi.org/10.1023/A:1018977225919>.

- 941 (36) Fini, A.; Fazio, G.; Gonzalez-Rodriguez, M.; Cavallari, C.; Passerini, N.; Rodriguez, L.
942 Formation of Ion-Pairs in Aqueous Solutions of Diclofenac Salts. *International Journal of*
943 *Pharmaceutics* **1999**, *187* (2), 163–173. [https://doi.org/10.1016/S0378-5173\(99\)00180-5](https://doi.org/10.1016/S0378-5173(99)00180-5).
- 944 (37) Sarveiya, V.; Templeton, J. F.; Benson, H. A. E. Ion-Pairs of Ibuprofen: Increased
945 Membrane Diffusion. *Journal of Pharmacy and Pharmacology* **2004**, *56* (6), 717–724.
946 <https://doi.org/10.1211/0022357023448>.
- 947 (38) Scherrer, R. A.; Donovan, S. F. Automated Potentiometric Titrations in KCl/ Water-
948 Saturated Octanol: Method for Quantifying Factors Influencing Ion-Pair Partitioning.
949 *Analytical Chemistry* **2009**, *81* (7), 2768–2778. <https://doi.org/10.1021/ac802729k>.
- 950 (39) Wenlock, M. C.; Potter, T.; Barton, P.; Austin, R. P. A Method for Measuring the
951 Lipophilicity of Compounds in Mixtures of 10. *SLAS Discovery* **2011**, *16* (3), 348–355.
952 <https://doi.org/10.1177/1087057110396372>.
- 953 (40) Fini, A.; Bassini, G.; Monastero, A.; Cavallari, C. Diclofenac Salts, VIII. Effect of the
954 Counterions on the Permeation through Porcine Membrane from Aqueous Saturated Solutions.
955 *Pharmaceutics* **2012**, *4* (3), 413–429. <https://doi.org/10.3390/pharmaceutics4030413>.
- 956 (41) Paternostre, M.; Meyer, O.; Grabielle-Madelmont, C.; Lesieur, S.; Ghanam, M.; Ollivon,
957 M. Partition Coefficient of a Surfactant between Aggregates and Solution: Application to the
958 Micelle-Vesicle Transition of Egg Phosphatidylcholine and Octyl Beta-D-Glucopyranoside.
959 *Biophysical Journal* **1995**, *69* (6), 2476–2488. [https://doi.org/10.1016/S0006-3495\(95\)80118-9](https://doi.org/10.1016/S0006-3495(95)80118-9).
960
- 961 (42) Pieńko, T.; Grudzień, M.; Taciak, P. P.; Mazurek, A. P. Cytisine Basicity, Solvation, LogP,
962 and LogD Theoretical Determination as Tool for Bioavailability Prediction. *Journal of*
963 *Molecular Graphics and Modelling* **2016**, *63*, 15–21.
964 <https://doi.org/10.1016/j.jmkgm.2015.11.003>.
- 965 (43) Quoc Hung, L. Electrochemical Properties of the Interface between Two Immiscible
966 Electrolyte Solutions: Part I. Equilibrium Situation and Galvani Potential Difference. *Journal*
967 *of Electroanalytical Chemistry and Interfacial Electrochemistry* **1980**, *115* (2), 159–174.
968 [https://doi.org/10.1016/S0022-0728\(80\)80323-8](https://doi.org/10.1016/S0022-0728(80)80323-8).
- 969 (44) Kakiuchi, T. Limiting Behavior in Equilibrium Partitioning of Ionic Components in
970 Liquid–Liquid Two-Phase Systems. *Anal. Chem.* **1996**, *68* (20), 3658–3664.
971 <https://doi.org/10.1021/ac960032y>.
- 972 (45) Disdier, Z.; Savoye, S.; Dagnelie, R. V. H. Effect of Solutes Structure and PH on the N-
973 Octanol/Water Partition Coefficient of Ionizable Organic Compounds. *Chemosphere* **2022**,
974 *304*. <https://doi.org/10.1016/j.chemosphere.2022.135155>.
- 975 (46) Berthod, A.; Carda-Broch, S.; Garcia-Alvarez-Coque, M. C. Hydrophobicity of Ionizable
976 Compounds. A Theoretical Study and Measurements of Diuretic Octanol-Water Partition
977 Coefficients by Countercurrent Chromatography. *Analytical Chemistry* **1999**, *71* (4), 879–888.
978 <https://doi.org/10.1021/AC9810563/ASSET/IMAGES/LARGE/AC9810563F00006.JPEG>.
- 979 (47) Âde, F.; Reymond, R.; Carrupt, P.-A.; Testa, B.; Girault, H. H. Charge and Delocalisation
980 Effects on the Lipophilicity of Protonable Drugs. *Chemistry – A European Journal* **1999**, *5*
981 (1), 39–48.

- 982 (48) Gobry, V.; Ulmeanu, S.; Reymond, F.; Bouchard, G.; Carrupt, P. A.; Testa, B.; Girault, H.
983 H. Generalization of Ionic Partition Diagrams to Lipophilic Compounds and to Biphasic
984 Systems with Variable Phase Volume Ratios. *Journal of the American Chemical Society* **2001**,
985 *123* (43), 10684–10690.
986 <https://doi.org/10.1021/JA015914F/ASSET/IMAGES/MEDIUM/JA015914FE00037.GIF>.
- 987 (49) Reymond, F.; Steyaert, G.; Carrupt, P. A.; Testa, B.; Girault, H. Ionic Partition Diagrams:
988 A Potential-PH Representation. *Journal of the American Chemical Society* **1996**, *118* (47),
989 11951–11957. https://doi.org/10.1021/JA962187T/SUPPL_FILE/JA11951.PDF.
- 990 (50) Cunha, R. D.; Ferreira, L. J.; Orestes, E.; Coutinho-Neto, M. D.; Almeida, J. M. de;
991 Carvalho, R. M.; Maciel, C. D.; Curutchet, C.; Homem-de-Mello, P. Naphthenic Acids
992 Aggregation: The Role of Salinity. *Computation* **2022**, *10* (10), 170.
993 <https://doi.org/10.3390/COMPUTATION10100170/S1>.
- 994 (51) Tshepelevitsh, S.; Hernits, K.; Leito, I. Prediction of Partition and Distribution Coefficients
995 in Various Solvent Pairs with COSMO-RS. *Journal of Computer-Aided Molecular Design*
996 **2018**, *32* (6), 711–722. <https://doi.org/10.1007/S10822-018-0125-Y/FIGURES/7>.
- 997 (52) Losada-Barreiro, S.; Paiva-Martins, F.; Bravo-Díaz, C. Partitioning of Antioxidants in
998 Edible Oil–Water Binary Systems and in Oil-in-Water Emulsions. *Antioxidants* **2023**, *12* (4),
999 828. <https://doi.org/10.3390/antiox12040828>.
- 1000 (53) Escher, B. I.; Abagyan, R.; Embry, M.; Klüver, N.; Redman, A. D.; Zarfl, C.; Parkerton,
1001 T. F. Recommendations for Improving Methods and Models for Aquatic Hazard Assessment
1002 of Ionizable Organic Chemicals. *Environmental Toxicology and Chemistry* **2020**, *39* (2), 269–
1003 286. <https://doi.org/10.1002/etc.4602>.
- 1004 (54) Hansima, M. A. C. K.; Zvomuya, F.; Amarakoon, I. Fate of Veterinary Antimicrobials in
1005 Canadian Prairie Soils – A Critical Review. *Science of The Total Environment* **2023**, *892*,
1006 164387. <https://doi.org/10.1016/j.scitotenv.2023.164387>.
- 1007 (55) Tsantili-Kakoulidou, A.; Panderi, I.; Csizmadia, F.; Darvas, F. Prediction of Distribution
1008 Coefficient from Structure. 2. Validation of Prolog D, an Expert System. *American*
1009 *Pharmaceutical Association* **1997**, *86* (10), 1173–1179.
- 1010 (56) Zamora, W. J.; Bertsch, E.; Suñer, S.; Pinheiro, S. Experimental N-Octanol/Water
1011 Partition/Distribution Coefficients Database for Small Molecules. **2023**.
1012 <https://doi.org/10.5281/ZENODO.7956685>.
- 1013 (57) Kim, S.; Chen, J.; Cheng, T.; Gindulyte, A.; He, J.; He, S.; Li, Q.; Shoemaker, B. A.;
1014 Thiessen, P. A.; Yu, B.; Zaslavsky, L.; Zhang, J.; Bolton, E. E. PubChem in 2021: New Data
1015 Content and Improved Web Interfaces. *Nucleic Acids Research* **2021**, *49* (D1), D1388–D1395.
1016 <https://doi.org/10.1093/NAR/GKAA971>.
- 1017 (58) Mishra, B.; Sankar, C.; Mishra, M. Polymer Based Solutions of Bupranolol Hydrochloride
1018 for Intranasal Systemic Delivery. *Journal of Drug Targeting* **2011**, *19* (3), 204–211.
1019 <https://doi.org/10.3109/1061186X.2010.492520>.
- 1020 (59) Bello, M. L.; Junior, A. M.; Freitas, C. A.; Moreira, M. L. A.; Costa, J. P. da; Souza, M.
1021 A. de; Santos, B. A. M. C.; Sousa, V. P. de; Castro, H. C.; Rodrigues, C. R.; Cabral, L. M.
1022 Development of Novel Montmorillonite-Based Sustained Release System for Oral

- 1023 Bromopride Delivery. *European Journal of Pharmaceutical Sciences* **2022**, *175*.
1024 <https://doi.org/10.1016/j.ejps.2022.106222>.
- 1025 (60) Bezençon, J.; Wittwer, M. B.; Cutting, B.; Smieško, M.; Wagner, B.; Kansy, M.; Ernst, B.
1026 PKa Determination by ¹H NMR Spectroscopy - An Old Methodology Revisited. *Journal of*
1027 *Pharmaceutical and Biomedical Analysis* **2014**, *93*, 147–155.
1028 <https://doi.org/10.1016/j.jpba.2013.12.014>.
- 1029 (61) Voigt, W.; Mannhold, R.; Limberg, J.; Blaschke, G. Interactions of Antiarrhythmics with
1030 Artificial Phospholipid Membranes. *Journal of Pharmaceutical Sciences* **1988**, *77* (12), 1018–
1031 1020.
- 1032 (62) Roseman, T. J.; Yalkowsky, S. H. Physicochemical Properties of Prostaglandin F_{2a}
1033 (Tromethamine Salt): Solubility Behavior, Surface Properties, and Ionization Constants.
1034 *Journal of Pharmaceutical Sciences* **1973**, *62* (10), 1680–1685.
- 1035 (63) Shalaeva, M.; Kenseth, J.; Lombardo, F.; Bastin, A. Measurement of Dissociation
1036 Constants (PK_a Values) of Organic Compounds by Multiplexed Capillary Electrophoresis
1037 Using Aqueous and Cosolvent Buffers. *Journal of Pharmaceutical Sciences* **2008**, *97* (7),
1038 2581–2606. <https://doi.org/10.1002/jps.21287>.
- 1039 (64) Qiang, Z.; Adams, C. Potentiometric Determination of Acid Dissociation Constants (PK_a)
1040 for Human and Veterinary Antibiotics. *Water Research* **2004**, *38* (12), 2874–2890.
1041 <https://doi.org/10.1016/j.watres.2004.03.017>.
- 1042 (65) Santos, T. de A. D. dos; Costa, D. O. da; Pita, S. S. da R.; Semaan, F. S. Potentiometric
1043 and Conductimetric Studies of Chemical Equilibria for Pyridoxine Hydrochloride in Aqueous
1044 Solutions: Simple Experimental Determination of PK_a Values and Analytical Applications to
1045 Pharmaceutical Analysis. *Ecl. Quím* **2010**, *35* (4), 81–86.
- 1046 (66) Morimoto, K.; Nagayasu, A.; Fukanoki, S.; Morisaka, K.; Hyon, S.-H.; Ikada, Y.
1047 Evaluation of Polyvinyl Alcohol Hydrogel as Sustained-Release Vehicle for Transdermal
1048 System of Bunitrolol-HCl-1. *DRUG DEVELOPMENT AND INDUSTRIAL PHARMACY*
1049 **1990**, *16* (1), 13–29.
- 1050 (67) Loftsson, T.; Thorisdóttir, S.; Fridriksdóttir, H.; Stefánsson, E. Enalaprilat and Enalapril
1051 Maleate Eyedrops Lower Intraocular Pressure in Rabbits. *Acta Ophthalmologica* **2010**, *88* (3),
1052 337–341. <https://doi.org/10.1111/j.1755-3768.2008.01495.x>.
- 1053 (68) Mannhold, R.; Dross, K. P.; Frekker, R.; Steen, van der. Drug Lipophilicity in QSAR
1054 Practice: I. A Comparison of Experimental with Calculative Approaches. *Quant. Struct. Act.*
1055 *Relat.* **1990**, *9*, 21–28.
- 1056 (69) Loftsson, T.; Vogensen, S. B.; Desbos, C.; Jansook, P. Carvedilol: Solubilization and
1057 Cyclodextrin Complexation: A Technical Note. *AAPS PharmSciTech* **2008**, *9* (2), 425–430.
1058 <https://doi.org/10.1208/s12249-008-9055-7>.
- 1059 (70) Kuntworbe, N.; Alany, R. G.; Brimble, M.; Al-Kassas, R. Determination of PK_a and
1060 Forced Degradation of the Indoloquinoline Antimalarial Compound Cryptolepine
1061 Hydrochloride. *Pharmaceutical Development and Technology* **2013**, *18* (4), 866–876.
1062 <https://doi.org/10.3109/10837450.2012.668554>.

- 1063 (71) Islam, M. S.; Narurkar, M. M. Solubility, Stability and Ionization Behaviour of Famotidine.
1064 *Journal of Pharmacy and Pharmacology* **1993**, *45* (8), 682–686.
1065 <https://doi.org/10.1111/j.2042-7158.1993.tb07088.x>.
- 1066 (72) Deng, Y.; Li, B.; Yu, K.; Zhang, T. Biotransformation and Adsorption of Pharmaceutical
1067 and Personal Care Products by Activated Sludge after Correcting Matrix Effects. *Science of*
1068 *the Total Environment* **2016**, *544*, 980–986. <https://doi.org/10.1016/j.scitotenv.2015.12.010>.
- 1069 (73) Franke, U.; Munk, A.; Wiese, M. Ionization Constants and Distribution Coefficients of
1070 Phenothiazines and Calcium Channel Antagonists Determined by a PH-Metric Method and
1071 Correlation with Calculated Partition Coefficients. *Journal of Pharmaceutical Sciences* **1999**,
1072 *88* (1), 89–95. <https://doi.org/10.1021/js980206m>.
- 1073 (74) Avdeef, A.; Box, K. J.; Comer, J. E. A.; Hibbert, C.; Tam, K. Y. PH-Metric LogP 10.
1074 Determination of Liposomal Membrane-Water Partition Coefficient of Ionizable Drugs.
1075 *Pharmaceutical Research* **1998**, *15* (2), 209–215.
- 1076 (75) Thanacoody, R. H. K. Thioridazine: The Good and the Bad. *Recent Patents on Anti-*
1077 *Infective Drug Discovery* **2011**, *6*, 92–98.
- 1078 (76) Martínez, V.; Maguregui, M. I.; Jiménez, R. M.; Alonso, R. M. Determination of the PK a
1079 Values of B-Blockers by Automated Potentiometric Titrations. *Journal of Pharmaceutical and*
1080 *Biomedical Analysis* **2000**, *23*, 459–468.
- 1081 (77) Huerta, B.; Jakimska, A.; Gros, M.; Rodríguez-Mozaz, S.; Barceló, D. Analysis of Multi-
1082 Class Pharmaceuticals in Fish Tissues by Ultra-High-Performance Liquid Chromatography
1083 Tandem Mass Spectrometry. *Journal of Chromatography A* **2013**, *1288*, 63–72.
1084 <https://doi.org/10.1016/j.chroma.2013.03.001>.
- 1085 (78) Fini, A.; Fazio, G.; Feroci, G. Solubility and Solubilization Properties of Non-Steroidal
1086 Anti-Inflammatory Drugs. *International Journal of Pharmaceutics* **1995**, *126* (1–2), 95–102.
1087 [https://doi.org/10.1016/0378-5173\(95\)04102-8](https://doi.org/10.1016/0378-5173(95)04102-8).
- 1088 (79) Jacka, M. R. *Clarke's Isolation and Identification of Drugs*, 2nd ed.; Moffat, A. C.,
1089 Jackson, J. V., Moss, M. S., Widdop, B., Greenfield, E. S., Eds.; Pharmaceutical Press, 2000.
- 1090 (80) Nakamura, Y.; Yamamoto, H.; Sekizawa, J.; Kondo, T.; Hirai, N.; Tatarazako, N. The
1091 Effects of PH on Fluoxetine in Japanese Medaka (*Oryzias Latipes*): Acute Toxicity in Fish
1092 Larvae and Bioaccumulation in Juvenile Fish. *Chemosphere* **2008**, *70* (5), 865–873.
1093 <https://doi.org/10.1016/j.chemosphere.2007.06.089>.
- 1094 (81) Schröder, W.; Andersson, J. T. Fast and Direct Method for Measuring 1-Octanol-Water
1095 Partition Coefficients Exemplified for Six Local Anesthetics. *Journal of Pharmaceutical*
1096 *Sciences* **2001**, *90* (12), 1948–1954. <https://doi.org/10.1002/JPS.1145>.
- 1097 (82) Avdeef, A. *Sirius Technical Application Notes (STAN)*; Sirius Analytical Instruments Ltd.,
1098 1994; Vol. 1.
- 1099 (83) Caron, G.; Steyaert, G.; Pagliara, A.; Âde, F.; Reymond, Â.; Crivori, P.; Gaillard, P.;
1100 Carrupt, P.-A.; Avdeef, A.; Comer, J.; Box, K. J.; Girault, H. H.; Testa, B. Structure-
1101 Lipophilicity Relationships of Neutral and Protonated b-Blockers Intra- and Intermolecular
1102 Effects in Isotropic Solvent Systems. [https://doi.org/10.1002/\(SICI\)1522-](https://doi.org/10.1002/(SICI)1522-2675(19990804)82:8)
1103 [2675\(19990804\)82:8](https://doi.org/10.1002/(SICI)1522-2675(19990804)82:8).

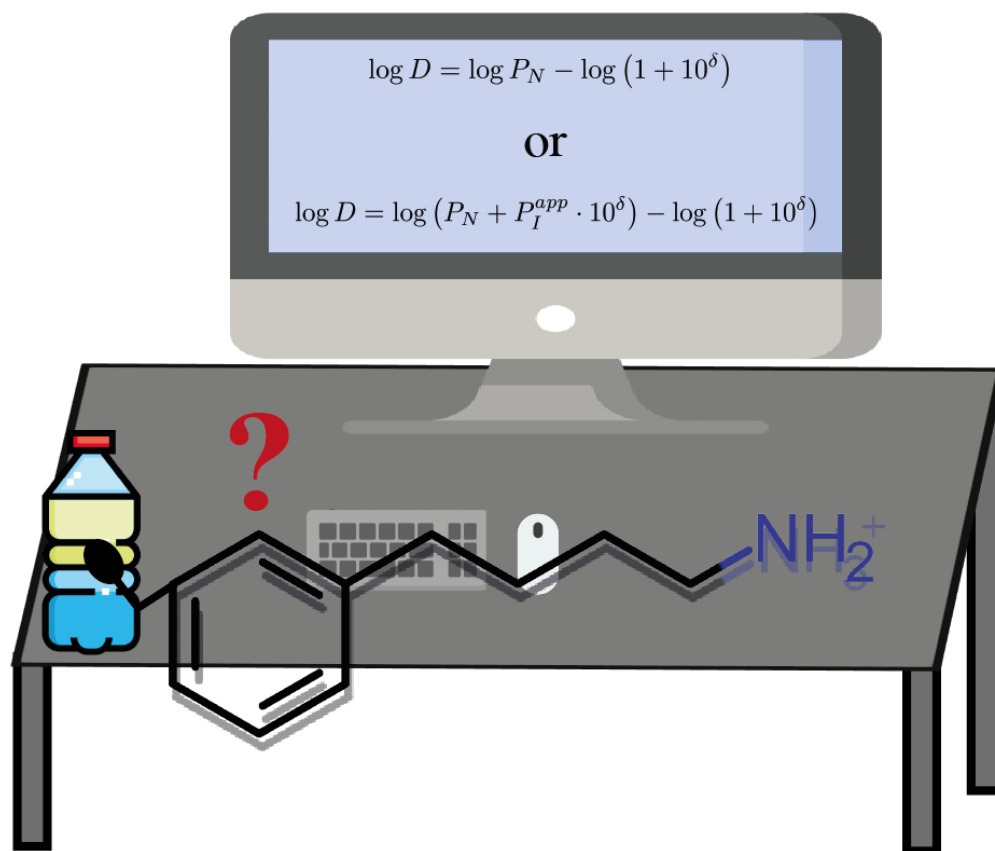
- 1104 (84) Avdeef, A. *Sirius Technical Application Notes (STAN)*; Sirius Analytical Instruments Ltd.,
1105 1995; Vol. 2.
- 1106 (85) Lombardo, F.; Shalaeva, M. Y.; Tupper, K. A.; Gao, F.; Abraham, M. H. ElogP(Oct): A
1107 Tool for Lipophilicity Determination in Drug Discovery. *Journal of Medicinal Chemistry*
1108 **2000**, *43* (15), 2922–2928.
1109 <https://doi.org/10.1021/JM0000822/ASSET/IMAGES/MEDIUM/JM0000822E00013.GIF>.
- 1110 (86) Winiwarter, S.; Bonham, N. M.; Ax, F.; Hallberg, A.; Lennernäs, H.; Karlén, A.
1111 Correlation of Human Jejunal Permeability (in Vivo) of Drugs with Experimentally and
1112 Theoretically Derived Parameters. A Multivariate Data Analysis Approach. *Journal of*
1113 *Medicinal Chemistry* **1998**, *41* (25), 4939–4949.
1114 https://doi.org/10.1021/JM9810102/SUPPL_FILE/JM9810102_S.PDF.
- 1115 (87) Slater, B.; McCormack, A.; Avdeef, A.; Comer, J. E. A. PH-Metric Log P. 4. Comparison
1116 of Partition Coefficients Determined by HPLC and Potentiometric Methods to Literature
1117 Values. *Journal of Pharmaceutical Sciences* **1994**, *83* (9), 1280–1283.
1118 <https://doi.org/10.1002/JPS.2600830918>.
- 1119 (88) Luger, P.; Daneck, K.; Engel, W.; Trummlitz, G.; Wagner, K. Structure and
1120 Physicochemical Properties of Meloxicam, a New NSAID. *European Journal of*
1121 *Pharmaceutical Sciences* **1996**, *4* (3), 175–187. [https://doi.org/10.1016/0928-0987\(95\)00046-](https://doi.org/10.1016/0928-0987(95)00046-1)
1122 [1](https://doi.org/10.1016/0928-0987(95)00046-1).
- 1123 (89) Takács-Novák, K.; Józán, M.; Hermeicz, I.; Szász, G. Lipophilicity of Antibacterial
1124 Fluoroquinolones. *International Journal of Pharmaceutics* **1992**, *79* (1–3), 89–96.
1125 [https://doi.org/10.1016/0378-5173\(92\)90099-N](https://doi.org/10.1016/0378-5173(92)90099-N).
- 1126 (90) Carda-Broch, S.; Berthod, A. PH Dependence of the Hydrophobicity of β -Blocker Amine
1127 Compounds Measured by Counter-Current Chromatography. *Journal of Chromatography A*
1128 **2003**, *995* (1–2), 55–66. [https://doi.org/10.1016/S0021-9673\(03\)00534-X](https://doi.org/10.1016/S0021-9673(03)00534-X).
- 1129 (91) Scott, D. Estimation of Distribution Coefficients from the Partition Coefficient and PKa.
1130 *Pharmaceutical Technology* **2002**, *26* (11).
- 1131 (92) Gulaboski, R.; Borges, F.; Pereira, C. M.; Natália, M.; Cordeiro, D. S.; Garrido, J.; Silva,
1132 A. F. Voltammetric Insights in the Transfer of Ionizable Drugs Across Biomimetic
1133 Membranes-Recent Achievements. *Combinatorial Chemistry & High Throughput Screening*
1134 **2007**, *10*, 514–526.
- 1135 (93) *pKa Plugin | Chemaxon Docs*. <https://docs.chemaxon.com/display/docs/pka-plugin.md>
1136 (accessed 2023-07-03).
- 1137 (94) *LogP and logD calculations | Chemaxon Docs*.
1138 <https://docs.chemaxon.com/display/docs/logp-and-logd-calculations.md> (accessed 2023-07-
1139 03).
- 1140 (122) Hall, L. H.; Kier, L. B. Electrotopological State Indices for Atom Types: A Novel
1141 Combination of Electronic, Topological, and Valence State Information. *Journal of Chemical*
1142 *Information and Computer Sciences* **1995**, *35* (6), 1039–1045.
1143 https://doi.org/10.1021/CI00028A014/ASSET/CI00028A014.FP.PNG_V03.

- 1144 (123) McLeod, A. I.; Xu, C.; Lai, Y. *Package “bestglm.”* [https://cran.r-](https://cran.r-project.org/web/packages/bestglm/bestglm.pdf)
1145 [project.org/web/packages/bestglm/bestglm.pdf](https://cran.r-project.org/web/packages/bestglm/bestglm.pdf).
- 1146 (124) Furnival, G. M.; Wilson, R. W. Regressions by Leaps and Bounds. *Technometrics* **1974**,
1147 *16* (4), 499–511. <https://doi.org/10.1080/00401706.1974.10489231>.
- 1148 (125) Burger, S. *Introduction to Machine Learning with R: Rigorous Mathematical Modeling*,
1149 1st ed.; O’Reilly, 2018.
- 1150 (126) Fauchere, J.; Pliska, V. Hydrophobic Parameters II of Amino Acid Side-Chains from the
1151 Partitioning of N-Acetyl-Amino Acid Amides. *Eur. J. Med. Chem.* **1983**, *18*.
- 1152 (127) Breiman, L. Random Forests. *Machine Learning* **2001**, *45* (1), 5–32.
1153 <https://doi.org/10.1023/A:1010933404324/METRICS>.
- 1154 (128) Breiman, L.; Cutler, A.; Liaw, A.; Wiener, M. *Package “randomForest.”* [https://cran.r-](https://cran.r-project.org/web/packages/randomForest/randomForest.pdf)
1155 [project.org/web/packages/randomForest/randomForest.pdf](https://cran.r-project.org/web/packages/randomForest/randomForest.pdf).
- 1156 (129) Cortes, C.; Vapnik, V.; Saitta, L. Support-Vector Networks. *Machine Learning 1995 20:3*
1157 **1995**, *20* (3), 273–297. <https://doi.org/10.1007/BF00994018>.
- 1158 (130) Boser, B. E.; Guyon, I. M.; Vapnik, V. N. Training Algorithm for Optimal Margin
1159 Classifiers. *Proceedings of the Fifth Annual ACM Workshop on Computational Learning*
1160 *Theory* **1992**, 144–152. <https://doi.org/10.1145/130385.130401>.
- 1161 (131) Meyer, D.; Dimitriadou, E.; Hornik, K.; Weingessel, A.; Leisch, F.; Chang, C.-C.; Lin, C.-
1162 C. *Package “e1071.”* <https://cran.r-project.org/web/packages/e1071/e1071.pdf>.
- 1163 (132) Xing, K.; Chatterjee, S.; Saito, T.; Gainaru, C.; Sokolov, A. P. Impact of Hydrogen
1164 Bonding on Dynamics of Hydroxyl-Terminated Polydimethylsiloxane. *Macromolecules* **2016**,
1165 *49* (8), 3138–3147. <https://doi.org/10.1021/acs.macromol.6b00262>.
- 1166 (133) Viayna, A.; Antermite, S. G.; de Candia, M.; Altomare, C. D.; Luque, F. J. Interplay
1167 between Ionization and Tautomerism in Bioactive β -Enamino Ester-Containing Cyclic
1168 Compounds: Study of Annulated 1,2,3,6-Tetrahydroazocine Derivatives. *J. Phys. Chem. B*
1169 **2020**, *124* (1), 28–37. <https://doi.org/10.1021/acs.jpcc.9b08904>.
- 1170 (134) Tielker, N.; Güssregen, S.; Kast, S. M. SAMPL7 Physical Property Prediction from EC-
1171 RISM Theory. *J Comput Aided Mol Des* **2021**, *35* (8), 933–941.
1172 <https://doi.org/10.1007/s10822-021-00410-9>.
- 1173 (135) Viayna, A.; Pinheiro, S.; Curutchet, C.; Luque, F. J.; Zamora, W. J. Prediction of N-
1174 Octanol/Water Partition Coefficients and Acidity Constants (PKa) in the SAMPL7 Blind
1175 Challenge with the IEFPCM-MST Model. *J Comput Aided Mol Des* **2021**, *35* (7), 803–811.
1176 <https://doi.org/10.1007/s10822-021-00394-6>.
- 1177 (136) Rodriguez, S. A.; Tran, J. V.; Sabatino, S. J.; Paluch, A. S. Predicting Octanol/Water
1178 Partition Coefficients and PKa for the SAMPL7 Challenge Using the SM12, SM8 and SMD
1179 Solvation Models. *J Comput Aided Mol Des* **2022**, *36* (9), 687–705.
1180 <https://doi.org/10.1007/s10822-022-00474-1>.
- 1181 (137) Wu, J.; Kang, Y.; Pan, P.; Hou, T. Machine Learning Methods for PKa Prediction of Small
1182 Molecules: Advances and Challenges. *Drug Discovery Today* **2022**, *27* (12), 103372.
1183 <https://doi.org/10.1016/j.drudis.2022.103372>.

- 1184 (138) Johnston, R. C.; Yao, K.; Kaplan, Z.; Chelliah, M.; Leswing, K.; Seekins, S.; Watts, S.;
1185 Calkins, D.; Chief Elk, J.; Jerome, S. V.; Repasky, M. P.; Shelley, J. C. Epik: PKa and
1186 Protonation State Prediction through Machine Learning. *J. Chem. Theory Comput.* **2023**, *19*
1187 (8), 2380–2388. <https://doi.org/10.1021/acs.jctc.3c00044>.
- 1188 (139) Šegatin, N.; Klofutar, C. Thermodynamics of the Solubility of Water in 1-Hexanol, 1-
1189 Octanol, 1-Decanol, and Cyclohexanol. *Monatshefte fur Chemie* **2004**, *135* (3), 241–248.
1190 <https://doi.org/10.1007/S00706-003-0053-X/METRICS>.
- 1191 (140) Reymond, F.; Chopineaux-Courtois, V.; Steyaert, G.; Bouchard, G.; Carrupt, P. A.; Testa,
1192 B.; Girault, H. H. Ionic Partition Diagrams of Ionisable Drugs: PH-Lipophilicity Profiles,
1193 Transfer Mechanisms and Charge Effects on Solvation. *Journal of Electroanalytical*
1194 *Chemistry* **1999**, *462* (2), 235–250. [https://doi.org/10.1016/S0022-0728\(98\)00418-5](https://doi.org/10.1016/S0022-0728(98)00418-5).
- 1195 (141) Burden, F. R. Molecular Identification Number for Substructure Searches. *Journal of*
1196 *Chemical Information and Computer Sciences* **1989**, *29* (3), 225–227.
1197 https://doi.org/10.1021/CI00063A011/ASSET/CI00063A011.FP.PNG_V03.
- 1198 (142) Burden, F. R. A Chemically Intuitive Molecular Index Based on the Eigenvalues of a
1199 Modified Adjacency Matrix. *Quant. Stmct-Act. Relat* **1997**, *16*, 3–314.
- 1200 (143) Pearlman, R. S.; Smith, K. M. Metric Validation and the Receptor-Relevant Subspace
1201 Concept. *Journal of Chemical Information and Computer Sciences* **1999**, *39* (1), 28–35.
1202 <https://doi.org/10.1021/CI980137X/ASSET/IMAGES/LARGE/CI980137XF00006.JPEG>.
- 1203 (144) Petitjean, M. Applications of the Radius-Diameter Diagram to the Classification of
1204 Topological and Geometrical Shapes of Chemical Compounds. *Journal of Chemical*
1205 *Information and Computer Sciences* **1992**, *32* (4), 331–337.
1206 https://doi.org/10.1021/CI00008A012/ASSET/CI00008A012.FP.PNG_V03.
- 1207
- 1208

1209 *TOC Graphics*

1210



1211

1212