

# Accuracy of macroscopic and microscopic pK<sub>a</sub> predictions of small molecules evaluated by the SAMPL6 Blind Challenge

Mehtap Işık (ORCID: [0000-0002-6789-952X](#))<sup>1,2\*</sup>, Ariën S. Rustenburg (ORCID: [0000-0002-3422-0613](#))<sup>1,3</sup>, Andrea Rizzi (ORCID: [0000-0001-7693-2013](#))<sup>1,4</sup>, M. R. Gunner (ORCID: [0000-0003-1120-5776](#))<sup>6</sup>, David L. Mobley (ORCID: [0000-0002-1083-5533](#))<sup>5</sup>, John D. Chodera (ORCID: [0000-0003-0542-119X](#))<sup>1</sup>

<sup>1</sup>Computational and Systems Biology Program, Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY 10065, United States; <sup>2</sup>Tri-Institutional PhD Program in Chemical Biology, Weill Cornell Graduate School of Medical Sciences, Cornell University, New York, NY 10065, United States; <sup>3</sup>Graduate Program in Physiology, Biophysics, and Systems Biology, Weill Cornell Medical College, New York, NY 10065, United States; <sup>4</sup>Tri-Institutional PhD Program in Computational Biology and Medicine, Weill Cornell Graduate School of Medical Sciences, Cornell University, New York, NY 10065, United States; <sup>5</sup>Department of Pharmaceutical Sciences and Department of Chemistry, University of California, Irvine, Irvine, California 92697, United States; <sup>6</sup>Department of Physics, City College of New York, New York NY 10031

\*For correspondence:  
[mehtap.isik@choderlab.org](mailto:mehtap.isik@choderlab.org) (MI)

## Abstract

The prediction of acid dissociation constants ( $pK_a$ ) is a prerequisite for predicting many other properties of a small molecule, such as its protein-ligand binding affinity, distribution coefficient ( $\log D$ ), membrane permeability, and solubility. The prediction of each of these properties requires knowledge of the relevant protonation states and solution free energy penalties of each state. The SAMPL6  $pK_a$  Challenge was the first time that a separate challenge was conducted for evaluating  $pK_a$  predictions as a part of the Statistical Assessment of Modeling of Proteins and Ligands (SAMPL). This challenge was motivated by the inaccuracies observed in prior physical property prediction challenges, such as SAMPL5  $\log D$  Challenge, caused by protonation state and  $pK_a$  prediction issues. The goal of the  $pK_a$  challenge was to assess the performance of contemporary  $pK_a$  prediction methods for drug-like molecules. The challenge set was composed of 24 small molecules that resembled fragments of kinase inhibitors, a number of which were multiprotic. Eleven research groups contributed blind predictions for a total of 37  $pK_a$  distinct prediction methods. In addition to blinded submissions, four widely used  $pK_a$  prediction methods were included in the analysis as reference methods. Collecting both microscopic and macroscopic  $pK_a$  predictions allowed in-depth evaluation of  $pK_a$  prediction performance. This article highlights deficiencies of typical  $pK_a$  prediction evaluation approaches when the distinction between microscopic and macroscopic  $pK_a$ s is ignored; in particular, we suggest more stringent evaluation criteria for microscopic and macroscopic  $pK_a$  predictions guided by the available experimental data. Top-performing submissions for macroscopic  $pK_a$  predictions achieved RMSE of 0.7–1.0  $pK_a$  units and included both quantum-mechanical and empirical approaches. The total number of extra or missing macroscopic  $pK_a$ s predicted by these submissions were fewer than 8 for 24 molecules. A large number of submissions had RMSE spanning 1–3  $pK_a$  units. Molecules with sulfur-containing heterocycles, iodo, and bromo groups suffered from less accurate  $pK_a$  predictions on average considering all methods evaluated. For a subset of molecules, we utilized experimentally determined microstates based on NMR to evaluate the dominant tautomer predictions for each macroscopic state. Prediction of dominant tautomers were a major source of error for microscopic  $pK_a$  predictions, especially errors in charged tautomers. The SAMPL6  $pK_a$  Challenge demonstrated the need for improving  $pK_a$  prediction methods for drug-like molecules, especially for challenging moieties and multiprotic molecules. The inaccuracy of  $pK_a$  predictions as observed in this challenge

41 can be detrimental to the performance of protein-ligand binding affinity predictions due to errors in predicted dominant charge  
42 and tautomeric states and errors in the calculation of free energy corrections for multiple protonation states of the ligand.

43

---

## 44 0.1 Keywords

45 SAMPL · blind prediction challenge · acid dissociation constant ·  $pK_a$  · small molecule · macroscopic  $pK_a$  · microscopic  $pK_a$  · macro-  
46 scopic protonation state · microscopic protonation state

## 47 0.2 Abbreviations

48 **SAMPL** Statistical Assessment of the Modeling of Proteins and Ligands

49  **$pK_a$**   $-\log_{10}$  acid dissociation equilibrium constant

50 **SEM** Standard error of the mean

51 **RMSE** Root mean squared error

52 **MAE** Mean absolute error

53  $\tau$  Kendall's rank correlation coefficient (Tau)

54 **R<sup>2</sup>** Coefficient of determination (R-Squared)

## 55 1 Introduction

56 The acid dissociation constant ( $K_a$ ) describes the protonation state equilibrium of a molecule given pH, and  $pK_a$  is the negative  
57 logarithmic form of  $K_a$ . Predicting  $pK_a$  is a prerequisite for predicting many other properties of small molecules such as their  
58 protein binding affinity, distribution coefficient ( $\log D$ ), membrane permeability, and solubility. Computer-aided drug design  
59 efforts help in the assessment of pharmaceutical and physicochemical properties assessment of virtual molecules to guide  
60 synthesis and prioritization decisions. Prior to synthesis of small molecule, experimental  $pK_a$  can not be measured. Therefore,  
61 accurate computational  $pK_a$  prediction methods are required.

62 Ionizable sites are found often in drug molecules and influence their pharmaceutical properties including target affinity,  
63 ADME/Tox, and formulation properties [1]. Drug molecules with titratable groups can exist in many different charge and proto-  
64 nation states based on the pH of the environment. Given that experimental data of protonation states and  $pK_a$  are often not  
65 available, we rely on predicted  $pK_a$  values to determine in which charge and protonation states the molecules exist and what  
66 are the relative populations of these states, so that we can assign the appropriate protonation state(s) in fixed-state calculations,  
67 or the appropriate solvent state weights/protonation penalty to calculations considering multiple states.

68 The pH of the human gut ranges between 1-8 and 74% of approved drugs can change ionization states within this physio-  
69 logical pH range [2]. Because of this,  $pK_a$  values of drug molecules provide essential information about their physicochemical  
70 and pharmaceutical properties. A wide distribution of acidic and basic  $pK_a$  values, ranging from 0 to 12, have been observed in  
71 approved drugs [1, 2].

72 Drug-like molecules present difficulties for  $pK_a$  prediction compared to simple monoprotic molecules. Drug-like molecules  
73 are frequently multiprotic, have large conjugated systems, heterocycles, and tautomerization. Besides, these larger molecules  
74 with conformational flexibility can have intramolecular hydrogen bonding which shifts  $pK_a$  values. Predicted shifts can be real or  
75 modeling artifacts due to collapsed conformations caused by deficiencies in solvation models. Yet predicting  $pK_a$ s of drug-like  
76 molecules accurately is a prerequisite for computational drug discovery and design.

77 Small molecule  $pK_a$  predictions influence computational protein-ligand binding affinities in multiple ways. Errors in  $pK_a$   
78 predictions can cause modeling the wrong charge and tautomerization states which affect hydrogen bonding opportunities and  
79 charge distribution of the ligand. The prediction of the dominant protonation state and relative population of minor states  
80 in the aqueous medium is dictated by the  $pK_a$  values. The relative free energy of different protonation states in the aqueous  
81 state is a function of pH, and it contributes to the overall protein-ligand affinity in the form of a free energy penalty of reaching  
82 higher energy protonation states [3]. Any error in predicting the free energy of a minor aqueous protonation state of a ligand  
83 that contributes to the complex formation will directly add to the error in the predicted binding free energy. Similarly for  $\log D$   
84 predictions, an inaccurate prediction of protonation states and their relative free energies will be detrimental to the accuracy of  
85 transfer free energy predictions.

86 For a monoprotic weak acid (HA) or base (B) dissociation equilibria shown in Equation 1, the acid dissociation constant is  
 87 expressed as in Equation 2, or, commonly, in its negative logarithmic form as in Equation 3. The ratio of ionization states can be  
 88 calculated with Henderson-Hasselbalch equations shown in Equation 4.



$$K_a = \frac{[A^-][H^+]}{[HA]} \quad K_b = \frac{[B][H^+]}{[B^+]} \quad (2)$$

$$pK_a = -\log_{10} K_a \quad (3)$$

$$pH = pK_a + \log_{10} \frac{[A^-]}{[HA]} \quad pH = pK_a + \log_{10} \frac{[B]}{[BH^+]} \quad (4)$$

89 The definition of  $pK_a$  diverges into two for multiprotic molecules: macroscopic  $pK_a$  and microscopic  $pK_a$  [4–6]. Macroscopic  
 90  $pK_a$  describes the equilibrium dissociation constant between different charged states of the molecule. Each charge state can be  
 91 composed of multiple tautomers. Macroscopic  $pK_a$  is about the deprotonation of the molecule, not the location of the titratable  
 92 group. A microscopic  $pK_a$  describes the acid dissociation equilibrium between individual tautomeric states of different charges.  
 93 (There is no  $pK_a$  defined between tautomers of the same charge as they have the same number of protons and their relative  
 94 populations are independent of pH). The microscopic  $pK_a$  determines the identity and distribution of tautomers within each  
 95 charge state. Thus, each macroscopic charge state of a molecule can be composed of multiple microscopic tautomeric states.  
 96 The microscopic  $pK_a$  value defined between two microstates captures the deprotonation of a single titratable group with a fixed  
 97 background protonation state of other titratable groups. In molecules with multiple titratable groups, the protonation state of  
 98 one group can affect the proton dissociation propensity of another functional group, therefore the same titratable group may  
 99 have different proton affinities (microscopic  $pK_a$  values) based on the protonation state of the rest of the molecule.

100 Different experimental methods can capture changes in the total charge or the location of individual protons, so they mea-  
 101 sure different definitions of  $pK_a$ s, as explained in more detail in prior work [7]. Most common  $pK_a$  measurement techniques  
 102 such as potentiometric and spectrophotometric methods measure macroscopic  $pK_a$ s while NMR measurements can determine  
 103 microscopic  $pK_a$ s by measuring microstate populations with respect to pH. Therefore, it is important to pay attention to the  
 104 source and definition of  $pK_a$  values to interpret their meaning correctly.

105 Many computational methods can predict both microscopic and macroscopic  $pK_a$ s. While experimental measurements more  
 106 often provide only macroscopic  $pK_a$ s, microscopic  $pK_a$  predictions are more informative for determining relevant microstates  
 107 (tautomers) of a molecule and their relative free energies. Predicted microstate populations can be converted to predicted  
 108 macroscopic  $pK_a$ s for direct comparison with experimentally obtained macroscopic  $pK_a$ s. In this paper, we explore approaches  
 109 to assess the performance of both macroscopic and microscopic  $pK_a$  predictions, taking advantage of available experimental  
 110 data.

111 Microscopic  $pK_a$  predictions can be converted to macroscopic  $pK_a$  predictions either directly with the equation 5 [8] or  
 112 through computing the macroscopic free energy of deprotonation between ionization states with charges N and N-1 via Boltz-  
 113 mann weighted sum of the relative free energy of microstates ( $G_i$ ) as in equations 6 and 7 [9].

$$K_a^{\text{macro}} = \sum_{j=1}^{N_{\text{deprot}}} \frac{1}{\sum_{i=1}^{N_{\text{prot}}} \frac{1}{K_{ij}^{\text{micro}}}} , \quad (5)$$

$$\Delta G_{N-1,N} = RT \ln \frac{\sum_i e^{-G_i/RT} \delta_{N_i, N-1}}{\sum_i e^{-G_i/RT} \delta_{N_i, N}} \quad (6)$$

$$pK_a = pH - \frac{\Delta G_{N-1,N}}{RT \ln 10} \quad (7)$$

114 In Equation 6  $\Delta G_{N-1,N}$  is the effective macroscopic protonation free energy.  $\delta_{N_i, N-1}$  is equal to 1 when the microstate i has a  
 115 total charge of N-1 and null otherwise.  $\delta_{N_i, N}$  is equal to 1 when the microstate i has a total charge of N and null otherwise. RT is  
 116 the ideal gas constant times the temperature.

## 117 1.1 Motivation for a blind $pK_a$ challenge

118 SAMPL (Statistical Assessment of the Modeling of Proteins and Ligands) is a series of annual computational prediction challenges  
119 for the computational chemistry community. The goal of SAMPL community is to evaluate the current performance of the models  
120 and to bring the attention of the quantitative biomolecular modeling field on problems that limit the accuracy of protein-ligand  
121 binding models.

122 SAMPL Challenges that focus on different physical properties so far have assessed intermolecular binding models of various  
123 protein-ligand and host-guest systems, as well as solvation models to predict hydration free energies and distribution coeffi-  
124 cients. Potential benefits of these challenges are motivating improvement computational methods and revealing unexpected  
125 error contributors by focusing on interesting test systems. SAMPL Challenges have demonstrated the effects of force field ac-  
126 curacy, sampling limitations, solvation modeling defects, and tautomer/protonation state predictions on protein-ligand binding  
127 predictions.

128 During the SAMPL5 log  $D$  Challenge, the performance of cyclohexane-water log  $D$  predictions was worse than expected  
129 and accuracy suffered when protonation states and tautomers were not taken into account [10, 11]. Many participants simply  
130 submitted log  $P$  predictions as if they were equal to log  $D$ , and many were not prepared to account for the contributions of  
131 different ionization states to distribution coefficient in their models. It was recognized that it does not matter how accurate  
132 the log  $P$  prediction is, if the protonation state effects are neglected. The calculations were improved by including free energy  
133 penalty of the neutral state which relies on obtaining an accurate  $pK_a$  prediction [? ]. With the motivation of deconvoluting the  
134 different sources of error contributing to the large errors observed in the SAMPL5 log  $D$  Challenge, we organized separate  $pK_a$   
135 and log  $P$  challenges in SAMPL6 [7, 12, 13]. For this iteration of the SAMPL challenge, we have taken one step back and isolated  
136 the problem of predicting aqueous protonation states.

137 This is the first time a blind  $pK_a$  prediction challenge has been fielded as part of SAMPL. In this challenge, we aimed to  
138 assess the performance of current  $pK_a$  prediction methods for drug-like molecules, investigate potential causes of inaccurate  
139  $pK_a$  estimates, and determine how much current level of accuracy might impact protein binding affinity predictions.

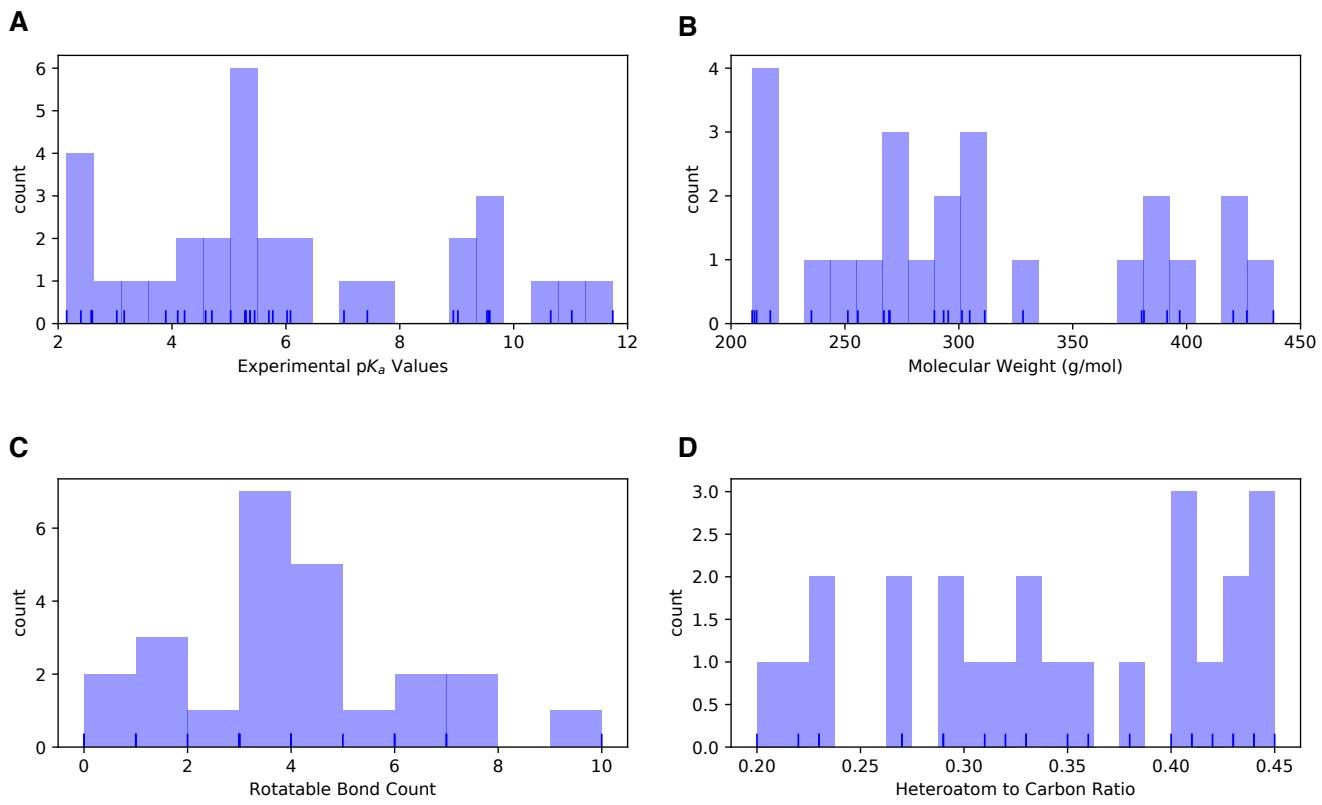
## 140 1.2 Approaches to predict small molecule $pK_a$ s

141 There is a large variety  $pK_a$  prediction methods developed for aqueous  $pK_a$  prediction of small molecules. Broadly we can di-  
142 vide  $pK_a$  predictions as knowledge-based empirical methods and physical methods. Empirical methods include the following  
143 categories: Database Lookup (DL) [14], Linear Free Energy Relationship (LFER) [15–17], Quantitative Structure-Property Rela-  
144 tionship (QSPR) [18–21], and Machine Learning approaches [22, 23]. DL methods rely on the principle that structurally similar  
145 compounds have similar  $pK_a$  values and utilize an experimental database of complete structures or fragments. The  $pK_a$  values  
146 of the most similar database entries are reported as the predicted  $pK_a$  of the query molecule. In the QSPR approach, the  $pK_a$   
147 values are predicted as a function of various quantitative molecular descriptors, and the parameters of the function are trained  
148 on experimental datasets. A function in the form of multiple linear regression is common, although more complex forms can  
149 also be used such as the artificial neural networks in ML methods. The LFER approach is the oldest  $pK_a$  prediction strategy. They  
150 use Hammett-Taft type equations to predict  $pK_a$  based on classification of the molecule to a parent class (associated with a base  
151  $pK_a$  value) and two parameters that describe how the base  $pK_a$  value must be modified given its substituents. Physical modeling  
152 of  $pK_a$  predictions require Quantum Mechanics (QM) models. QM methods are often utilized together with linear empirical cor-  
153 rections (LEC) that are designed to rescale and unbias QM predictions for better accuracy. Classical molecular mechanics-based  
154  $pK_a$  prediction methods are not feasible as deprotonation is a covalent bond breaking event that can only be captured by QM.  
155 Constant-pH molecular dynamics methods can calculate  $pK_a$  shifts in large biomolecular systems where there is low degree of  
156 coupling between protonation sites and linear summation of protonation energies can be assumed [? ]. However, this approach  
157 can not be applied to small organic molecule due to the high degree of coupling between protonation sites.

## 158 2 Methods

### 159 2.1 Design and logistics of the SAMPL6 $pK_a$ Challenge

160 The SAMPL6  $pK_a$  Challenge was conducted as a blind prediction challenge and focused on predicting aqueous  $pK_a$  values of 24  
161 small molecules. The challenge set was composed of molecules that resemble fragments of kinase inhibitors. Heterocycles that  
162 are frequently found in FDA-approved kinase inhibitors were represented in this set. The compound selection process was  
163 described in depth in the prior publication reporting SAMPL6  $pK_a$  Challenge experimental data collection [7]. The distribution of  
164 molecular weights, experimental  $pK_a$  values, number of rotatable bonds, and heteroatom to carbon ratio are depicted in Fig. 1.



**Figure 1. Distribution of molecular properties of 24 compounds in SAMPL6  $pK_a$  Challenge.** **A** Histogram of spectrophotometric  $pK_a$  measurements collected with Sirius T3 [7]. The overlayed carpet plot indicates the actual values. Five compounds have multiple measured  $pK_a$ s in the range of 2-12. **B** Histogram of molecular weights calculated for the neutral state of the compounds in SAMPL6 set. Molecular weights were calculated by neglecting counter ions. **C** Histogram of the number of non-terminal rotatable bonds in each molecule. **D** The histogram of the ratio of heteroatom (non-carbon heavy atoms including, O, N, F, S, Cl, Br, I) count to the number of carbon atoms.

165 The challenge molecule set was composed of 17 small molecules with limited flexibility (less than 5 non-terminal rotatable bonds)  
166 and 7 molecules with 5-10 non-terminal rotatable bonds. The distribution of experimental  $pK_a$  values ranged between 2-12 and  
167 roughly uniform. 2D representations of all compounds were provided in Fig. 5. Drug-like molecules are often larger and more  
168 complex than the ones used in this study. We limited the size and the number of rotatable bonds of compounds to create  
169 molecule set of intermediate difficulty.

170 The dataset composition and details of the  $pK_a$  measurement technique without the identity of the small molecules, were  
171 announced about a month before the challenge start time. Experimental macroscopic  $pK_a$  measurements were collected with  
172 the spectrophotometric method of Sirius T3, at room temperature in ionic strength-adjusted water with 0.15 M KCl [7]. The  
173 instructions for participation and the identity of the challenge molecules were released at the challenge start date (October 25,  
174 2017). A table of molecule IDs (in the form of SM##) and their canonical isomeric SMILES was provided as input. Blind prediction  
175 submissions were accepted until January 22, 2018.

176 Following the conclusion of the blind challenge, the experimental data was made public on January 23, 2018. The SAMPL  
177 organizers and participants gathered at the Second Joint D3R/SAMPL Workshop, at UC San Diego, La Jolla, CA on February 22-23,  
178 2018 to share results. The workshop aimed to create an opportunity for participants to have discussions, evaluate the results  
179 and lessons of the challenge together. The participants reported their results and their own evaluations in a special issue of the  
180 Journal of Computer-Aided Molecular Design [24].

181 While designing this first  $pK_a$  prediction challenge, we did not know the optimal format to capture  $pK_a$  predictions of par-  
182 ticipants. We wanted capture all necessary information that will aid the evaluation of  $pK_a$  predictions at the submission stage.  
183 Our strategy was to directly evaluate macroscopic  $pK_a$  predictions comparing them to experimental macroscopic  $pK_a$  values and  
184 to use collected microscopic  $pK_a$  prediction data for more in-depth diagnostics of method performance. Therefore, we asked  
185 participants to submit their predictions in three different submission types:

- 186 • **Type I:** microscopic  $pK_a$  values and related microstate pairs
- 187 • **Type II:** fractional microstate populations as a function of pH in 0.1 pH increments
- 188 • **Type III:** macroscopic  $pK_a$  values

189 For each submission type, a machine-readable submission file template was specified. For type I submissions, participants  
190 were asked to report microstate ID of protonated state, microstate ID of deprotonated state, microscopic  $pK_a$ , and microscopic  
191  $pK_a$  SEM. The reason and method of microstate enumeration are discussed further in Section 2.2 "Enumeration of Microstates".  
192 The SEM captures the statistical uncertainty of the prediction method. Microstate IDs were preassigned identifiers for each  
193 microstate in the form of SM##\_micro##. For type II submissions, the submission format included a table that started with a  
194 microstate ID column and a set of columns reporting the natural logarithm of fractional microstate population values of each  
195 predicted microstate for 0.1 pH increments between pH 2 and 12. For type III submissions participants were asked to report  
196 molecule ID, macroscopic  $pK_a$ , macroscopic  $pK_a$  SEM.

197 It was mandatory to submit predictions for all fields for each prediction, but it was not mandatory to submit predictions for  
198 all the molecules or all the submission types. Although we accepted submissions with partial sets of molecules, it would have  
199 been a better choice to require predictions for all the molecules for a better comparison of overall method performance. The  
200 submission files also included fields for naming the method, listing the software utilized, and a free text method section for the  
201 detailed documentation of each method.

202 Participants were allowed to submit predictions with multiple methods as long as they created separate submission files.  
203 Anonymous participation was allowed in the challenge, however, all participants opted to make their submissions public. All  
204 blind submissions were assigned a unique 5-digit alphanumeric submission ID, which will be used throughout this paper. Unique  
205 IDs were also assigned when multiple submissions exist for different submissions types of the same method such as microscopic  
206  $pK_a$ (type I) and macroscopic  $pK_a$  (type III). These submission IDs were also reported in the evaluation papers of participants and  
207 allow cross-referencing. Submission IDs, participant provided method names, and method categories are presented in Table 1.  
208 In many cases, multiple types of submissions of the same method were provided by participants as challenge instructions re-  
209 quested. Although each prediction set was assigned a separate submission ID, we have matched the submissions that originated  
210 from the same method according to the reports of the participants, for cases where multiple sets of predictions came from a  
211 given method. Submission IDs for both macroscopic (type III) and microscopic (type I)  $pK_a$  predictions for each method are  
212 shown in Table 1.

## 213 2.2 Enumeration of microstates

214 To capture both the  $pK_a$  value and titration position of microscopic  $pK_a$  predictions, we needed microscopic  $pK_a$  values to be  
215 reported together with a pair of microstates which describe the protonated and deprotonated states of each microscopic transi-  
216 tion. String representations of molecules such as canonical SMILES with explicit hydrogens can be written, however, there can be  
217 inconsistencies between the interpretation of canonical SMILES written by different software and algorithms. To avoid complica-  
218 tions while reading microstate structure files from different sources, we have decided that the safest route was pre-enumerating  
219 all possible microstates of challenge compounds, assigning the microstates IDs to each in the form of SM##\_micro##, and re-  
220 quire participants to report microscopic  $pK_a$  values along with microstate pairs specified by the provided microstates IDs.

221 We created initial sets of microstates with Epik [25] and OpenEye QUACPAC [26] and took the union of results. Microstates  
222 with Epik were generated using Schrödinger Suite v2016-4, running Epik to enumerate all tautomers within 20  $pK_a$  units of pH 7.  
223 For enumerating microstates with OpenEye QUACPAC, we had to first enumerate formal charges and for each charge enumerate  
224 all possible tautomers using the settings of maximum tautomer count 200, level 5, with carbonyl hybridization set to False. Then  
225 we created a union of all enumerated states written as canonical isomeric SMILES. Even though resonance structures correspond  
226 to different canonical isomeric SMILES they are not different microstates, therefore it was necessary to remove resonance struc-  
227 tures that were replicates of the same tautomer. To detect resonance structures we converted canonical isomeric SMILES to  
228 InChI hashes with explicit and fixed hydrogen layer. Structures that describe the same tautomer but different resonance states  
229 lead to explicit hydrogen InChI hashes that are identical, allowing replicates to be removed. The Jupyter Notebook used for the  
230 enumeration of microstates is provided in supplementary information. We provided microstate ID tables with canonical SMILES  
231 and 2D-depictions to aid participants in matching predicted structures to microstate IDs. Canonical SMILES representation was  
232 selected over canonical isomeric SMILES, because resonance and geometric isomerism do not lead to different microstates ac-  
233 cording to our working microstate definition. The only exception was for molecules SM20 which should be consistently modeled  
234 as the E-isomer.

235 In SAMPL6 Challenge, participants came up with new microstates that were not present in the initial list that we provided. De-  
236 spite pooling together enumerated charge states and tautomers generated by both Epik and OpenEye QUACPAC to our surprise  
237 the microstate lists were still incomplete. Based on participant requests for new microstates, we iteratively had to update the list  
238 of microstates and assign new microstate IDs. Every time we received a request, we shared the updated microstate ID lists with  
239 all the challenge participants. Some participants updated their  $pK_a$  prediction by including the newly added microstates in their  
240 calculations. In the future, developing a better algorithm that can enumerate all possible microstates (not just the ones with  
241 significant populations) would be very beneficial for anticipating microstates that may be predicted by  $pK_a$  prediction methods.

242 A microscopic  $pK_a$  definition was provided in challenge instructions for clarity as follows: Physically meaningful microscopic  
243  $pK_a$ s are defined between microstate pairs that can interconvert by single protonation/deprotonation event of only one titrable  
244 group. So, microstate pairs should have total charge (absolute) difference of 1 and only one heavy atom that differs in the number  
245 of bound hydrogens, regardless of resonance state or geometric isomerism. All geometric isomer and resonance structure  
246 pairs that have the same number of hydrogens bound to equivalent heavy atoms are grouped in the same microstate. Pairs of  
247 resonance structures and geometric isomers (cis/trans, stereo) are not considered as different microstates, as long as there is  
248 no change in the number of hydrogens bound to each heavy atom. Transitions where there are shifts in the position of protons  
249 coupled to changes in the number of protons were also not considered as microscopic  $pK_a$ s [27]. Since we wanted participants  
250 to report only microscopic  $pK_a$ s that describe single deprotonation events (in contrast to transitions between microstates that  
251 are different in terms of two or more titratable protons), we have also provided a pre-enumerated list of allowed microstate  
252 pairs.

253 Provided microstate ID and microstate pair lists were intended to be used for reporting microstate IDs and to aid parsing  
254 of submissions. The enumerated lists of microstates were not created with the intent to guide computational predictions. This  
255 was clearly stated in the challenge instructions. However, we noticed that some participants still used the microstate lists as  
256 an input for their  $pK_a$  predictions as we received complaints from participants that due to our updates to microstate lists they  
257 needed to repeat their calculations. This would not have been an issue if participants used  $pK_a$  prediction protocols that did not  
258 rely on an external pre-enumerated list of microstates as an input. None of the participants have reported this dependency in  
259 their method descriptions explicitly, therefore it was also not obvious how participants were using the provided states in their  
260 predictions. We could not identify which submissions used these enumerated microstate lists as input for predictions and which  
261 have followed the challenge instructions and relied only on their prediction method to generate microstates.

## 2.3 Evaluation approaches

Since the experimental data for the challenge was mainly composed of macroscopic  $pK_a$  values of both monoprotic and multiprotoic compounds, evaluation of macroscopic and microscopic  $pK_a$  predictions was not straightforward. For a subset of 8 molecules, the dominant microstate sequence could be inferred from NMR. For the rest of the molecules, the only experimental information available was the macroscopic  $pK_a$  value. The experimental data in the form of macroscopic  $pK_a$ s did not provide any information on which group(s) are being titrated, the microscopic  $pK_a$  values, the identity of the associated macrostates (which total charge) or microstates (which tautomers). Also, experimental data did not provide any information about the charge state of protonated and deprotonated species associated with each macroscopic  $pK_a$ . Typically charges of states associated with experimental  $pK_a$  values are assigned based on  $pK_a$  predictions, not experimental evidence, but we did not utilize such computational charge assignment. For a fair performance comparison between methods, we avoided relying on any particular  $pK_a$  prediction to assist the interpretation of the experimental reference data. This choice complicated the  $pK_a$  prediction analysis, especially regarding how to pair experimental and predicted  $pK_a$ s for error analysis. We adopted various evaluation strategies guided by the experimental data. To compare macroscopic  $pK_a$  predictions to experimental values we had to utilize numerical matching algorithms before we could calculate performance statistics. For the subset of molecules with experimental data about microstates, we used microstate based matching. These matching methods were described further in the next section.

Three types of submissions were collected during the SAMPL6  $pK_a$  Challenge. We have only utilized the type I (microscopic  $pK_a$  value and microstate IDs) and the type III (macroscopic  $pK_a$  value) predictions in this article. Type I submissions contained the same prediction information as the type II submissions which reported the fractional population of microstates with respect to pH. We collected type II submissions in order to capture relative populations of microstates, not realizing they were redundant. The microscopic  $pK_a$  predictions collected in type I submissions capture all the information necessary to calculate type II submissions. Therefore, we did not use type II submissions for challenge evaluation. In theory, type III (macroscopic  $pK_a$ ) predictions can also be calculated from type I submissions. But collecting type III submissions allowed participation of  $pK_a$  prediction methods that directly predict macroscopic  $pK_a$ s without considering microspeciation and methods that apply special empirical corrections for macroscopic  $pK_a$  predictions.

### 2.3.1 Matching algorithms for pairing predicted and experimental $pK_a$ s

Macroscopic  $pK_a$  predictions can be calculated from microscopic  $pK_a$ s for direct comparison to experimental macroscopic  $pK_a$  values. One major question must be answered to allow this comparison: How should we match predicted macroscopic  $pK_a$ s to experimental macroscopic  $pK_a$ s when there could multiple  $pK_a$  values reported for a given molecule? For example, experiments on SM18 showed three macroscopic  $pK_a$ s, but prediction of *xvxzd* method reported two macroscopic  $pK_a$  values. There were also examples of the opposite situation with more predicted  $pK_a$  values than experimentally determined macroscopic  $pK_a$ s: One experimental  $pK_a$  was measured for SM02, but two macroscopic  $pK_a$ s were predicted by *xvxzd* method. The experimental and predicted values must be paired before any prediction error can be calculated, even though there was not any experimental information regarding underlying tautomer and charge states.

Knowing the charges of macrostates would have guided the pairing between experimental and predicted macroscopic  $pK_a$ s, however, not all experimental  $pK_a$  measurements can determine determine the charge of protonation states. The potentiometric  $pK_a$  measurements just captures the relative charge change between macrostates, but not the absolute value of the charge. Thus, our experimental data did not provide any information that would indicate the titration site, the overall charge, or the tautomer composition of macrostate pairs that are associated with each measured macroscopic  $pK_a$  that can guide the matching between predicted and experimental  $pK_a$  values.

For evaluating macroscopic  $pK_a$  predictions taking the experimental data as reference, Fraczkiewicz et al. delineated recommendations for fair comparative analysis of computational  $pK_a$  predictions [22]. They recommended that, in the absence of any experimental information that would aid in matching, experimental and computational  $pK_a$ s should be matched preserving the order of  $pK_a$  values and minimizing the sum of absolute errors.

We picked the Hungarian matching algorithm [28, 29] to match experimental and predicted macroscopic  $pK_a$ s with a squared error cost function as suggested by Kiril Lanevskij via personal communication. The algorithm is available in the SciPy package (`scipy.optimize.linear_sum_assignment`) [30]. This matching algorithm provides optimum global assignment that minimizes the linear sum of squared errors of all pairwise matches. We selected the squared error cost function instead of the absolute error cost function to avoid misordered matches. For instance, for a molecule with experimental  $pK_a$  values of 4 and 6, and predicted  $pK_a$ s of 7 and 8, Hungarian matching with absolute error cost function would match 6 to 7 and 4 to 9. Hungarian matching with squared error cost would match 4 to 7 and 6 to 9, preserving the increasing  $pK_a$  value order between experimental and

predicted values. A weakness of this approach would be failing to match the experimental value of 6 to predicted value of 7 if that was the correct match based on underlying macrostates. But the underlying pair of states were unknown to us both because the experimental data did not determine which charge states the transitions were happening between and also because we did not collect the pair of macrostates associated with each  $pK_a$  predictions in submissions. Requiring this information for macroscopic  $pK_a$  predictions in future SAMPL challenges would allow for better comparison between predictions, even if experimental assignment of charges is not possible. There is no perfect solution to the numerical  $pK_a$  assignment problem, but we tried to determine the fairest way to penalize predictions based on their numerical deviation from the experimental values.

For the analysis of microscopic  $pK_a$  predictions we adopted a different matching approach. For the eight molecules for which we had the requisite data for this analysis, we utilized the dominant microstate sequence inferred from NMR experiments to match computational predictions and experimental  $pK_a$ s. We will refer to this assignment method as microstate matching, where the experimental  $pK_a$  value is matched to the computational microscopic  $pK_a$  value which was reported for the dominant microstate pair observed for each transition. We have compared the results of Hungarian matching and microstate matching.

Inevitably the choice of matching algorithms to assign experimental and predicted values has an impact on the calculation of performance statistics. We believe the Hungarian algorithm for numerical matching of unassigned  $pK_a$  values and microstate-based matching when experimental microstates are known were the best choices, providing the most unbiased matching without introducing assumptions outside of the experimental data.

### 2.3.2 Statistical metrics for submission performance

A variety of accuracy and correlation statistics were considered for analyzing and comparing the performance of prediction methods submitted to the SAMPL6  $pK_a$  Challenge. Calculated performance statistics of predictions were provided to participants before the workshop. Details of the analysis and scripts are maintained on the SAMPL6 Github Repository (described in Section 5).

There are six error metrics reported for the numerical error of the  $pK_a$  values: the root-mean-squared error (RMSE), mean absolute error (MAE), mean error (ME), coefficient of determination ( $R^2$ ), linear regression slope (m), and Kendall's Rank Correlation Coefficient ( $\tau$ ). Uncertainty in each performance statistic was calculated as 95% confidence intervals estimated by bootstrapping over predictions with 10000 bootstrap samples. Calculated errors statistics of all methods can be found in Table S2 for macroscopic  $pK_a$  predictions and Tables S4 and S4 for microscopic  $pK_a$  predictions.

In addition to the numerical error aspect of the  $pK_a$  values, we also evaluated predictions in terms of their ability to capture the correct macrostates (ionization states) and microstates (tautomers of each ionization state) to the extent possible from the available experimental data. For macroscopic  $pK_a$ s, experiments did not provide any evidence of the identity of the ionization states. However, the number of ionization states indicates the number of macroscopic  $pK_a$ s that exists between the experimental range of 2.0-12.0. For instance, SM14 has two experimental  $pK_a$ s and therefore three different charge states observed between pH 2.0 and 12.0. If a prediction reported 4 macroscopic  $pK_a$ s, it is clear that this method predicted an extra ionization state. With this perspective we reported the number of unmatched experimental  $pK_a$ s (the number of missing  $pK_a$  predictions, i.e. missing ionization states) and the number of unmatched predicted  $pK_a$ s (the number of extra  $pK_a$  predictions, i.e. extra ionization states) after Hungarian matching. The later count was restricted to only predictions with  $pK_a$  values between 2 and 12 because that was the range of the experimental method. Errors in extra or missing  $pK_a$  prediction errors highlight failure to predict the correct number of ionization states within a pH range.

For the evaluation of microscopic  $pK_a$  predictions, taking advantage of the available dominant microstate sequence data for a subset of 8 compounds, we calculated the dominant microstate prediction accuracy which is the ratio of correct dominant tautomer predictions for each charge state divided by the total number of dominant tautomer predictions. Dominant microstate accuracy was calculated over all experimentally detected ionization states of each molecule which were part of this analysis. In order to extract the sequence of dominant microstates from the microscopic  $pK_a$  predictions sets, we calculated the relative free energy of microstates selecting a neutral tautomer and pH 0 as reference following Equation 8. Calculation of relative free energy of microstates was explained in more detail in a previous publication [27].

The relative free energy of a state with respect to reference state B at pH 0.0 (arbitrary pH value selected as reference) can be calculated as follows:

$$\Delta G_{AB} = \Delta m_{AB} RT \ln 10 (pH - pK_a) \quad (8)$$

$\Delta m_{AB}$  is equal to the number protons in state A minus that in state B. R and T indicate the molar gas constant and temperature, respectively. By calculating relative free energies of all predicted microstates with respect to the same reference state and pH,

360 we were able to determine the sequence of predicted dominant microstates. The dominant tautomer of each charge state  
361 was determined as the microstate with the lowest free energy in the subset of predicted microstates of each ionization state.  
362 This approach is feasible because the relative free energy of tautomers of the same ionization state is independent of pH and  
363 therefore the choice of reference pH is arbitrary.

364 We created a shortlist of top-performing methods for macroscopic and microscopic  $pK_a$  predictions. The top macroscopic  $pK_a$   
365 predictions were selected if they ranked in the top 10 consistently according to two error metrics (RMSE, MAE) and two correlation  
366 metrics (R-Squared, and Kendall's Tau), while also having fewer than eight missing or extra macroscopic  $pK_a$ s for the entire  
367 molecule set (eight macrostate errors correspond to macrostate prediction mistake in roughly one third of the 24 compounds).  
368 These methods are presented in Table 2. A separate list of top-performing methods was constructed for microscopic  $pK_a$  with  
369 the following criteria: ranking in the top 10 methods when ranked by accuracy statistics (RMSE and MAE) and perfect dominant  
370 microstate prediction accuracy. These methods are presented in Table 3.

371 In addition to comparing the performance of methods, we also wanted to compare  $pK_a$  prediction performance for each  
372 molecule to determine which molecules were the most challenging for  $pK_a$  predictions considering all the methods in the chal-  
373 lenge. For this purpose, we plotted prediction error distributions of each molecule calculated over all prediction methods. We  
374 also calculated MAE for each molecule over all prediction sets as well as for predictions from each method category separately.

## 375 2.4 Reference calculations

376 Including a null model is helpful in comparative performance analysis of predictive methods to establish what the performance  
377 statistics look like for a baseline method for the specific dataset. Null models or null predictions employ a simple prediction  
378 model which is not expected to be particularly successful, but it provides a simple point of comparison for more sophisticated  
379 methods. The expectation or goal is for more sophisticated or costly prediction methods to outperform the predictions from a  
380 null model, otherwise the simpler null model would be preferable. In SAMPL6  $pK_a$  Challenge there were two blind submissions  
381 using database lookup methods that were submitted to serve as null predictions. These methods, with submission IDs 5nm4j and  
382 5nm4j both used OpenEye pKa-Prospector database to find the most similar molecule to query molecule and simply reported its  
383  $pK_a$  as the predicted value. Database lookup methods with a rich experimental database do present a challenging null model to  
384 beat, however, due to the accuracy level needed from  $pK_a$  predictions for computer-aided drug design we believe such methods  
385 provide an appropriate performance baseline that physical and empirical  $pK_a$  prediction methods should strive to outperform.

386 We also included additional reference calculations in the comparative analysis to provide more perspective. Some widely  
387 used methods by academia and industry were missing from the blind challenge submission. Therefore, we included those meth-  
388 ods as reference calculations: Schrödinger/Epik (nb007, nb008, nb010), Schrödinger/Jaguar (nb011, nb013), Chemaxon/Chemicalize  
389 (nb015), and Molecular Discovery/MoKa (nb016, nb017). Epik and Jaguar  $pK_a$  predictions were collected by Bas Rustenburg, Chem-  
390 icalize predictions by Mehtap Isik, and MoKa predictions by Thomas Fox. All were done after the challenge deadline avoiding  
391 any alterations to their respective standard procedures and any guidance from experimental data. Experimental data was pub-  
392 licly available before these calculations were complete, therefore reference calculations were not formally considered as blind  
393 submissions.

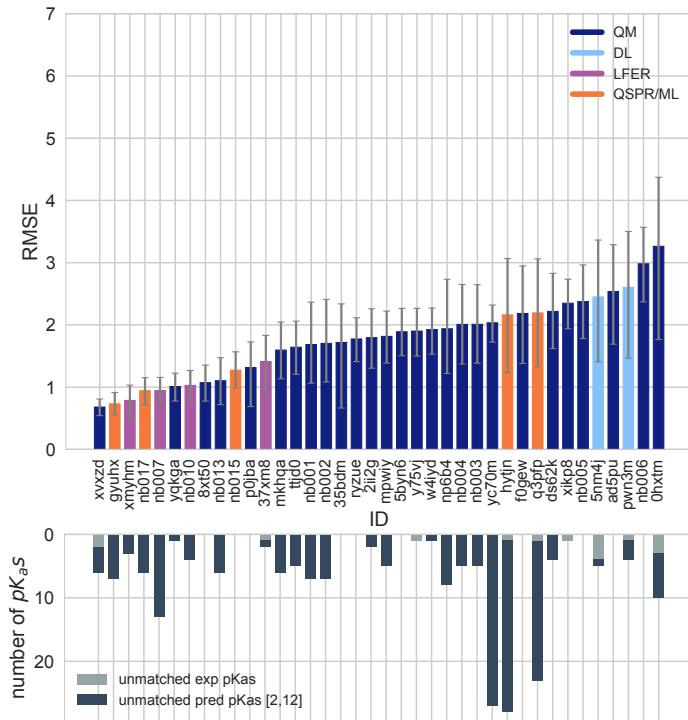
394 All figures and statistics tables in this manuscript include reference calculations. As the reference calculations were not formal  
395 submissions, these were omitted from formal ranking in the challenge, but we present plots in this article which show them for  
396 easy comparison. These are labeled with submission IDs of the form nb### to clearly indicate non-blind reference calculations.

## 397 3 Results and Discussion

398 Participation in SAMPL6  $pK_a$  Challenge was high with 11 research groups contributing  $pK_a$  prediction sets for 37 methods. A  
399 large variety of  $pK_a$  prediction methods were represented in the SAMPL6 Challenge. We categorized these submissions into  
400 four method classes: database lookup (DL), linear free energy relationship (LFER), quantitative structure-property relationship  
401 or machine learning (QSPR/ML), and quantum mechanics (QM). Quantum mechanics models were subcategorized into QM  
402 methods with and without linear empirical correction (LEC), and combined quantum mechanics and molecular mechanics (QM  
403 + MM). Table 1 presents method names, submission IDs, method categories, and also references for each approach. Integral  
404 equation-based approaches (e.g. EC-RISM) were also evaluated under the Physical (QM) category. There were 2 DL, 4 LFER, and  
405 5 QSPR/ML methods represented in the challenge, including the reference calculations. The majority of QM calculations include  
406 linear empirical corrections (22 methods in QM + LEC category), and only 5 QM methods were submitted without any empirical  
407 corrections. There were 4 methods that used a mixed physical modeling approach of QM + MM.

The following sections present a detailed performance evaluation of blind submissions and reference prediction methods for macroscopic and microscopic  $pK_a$  predictions. Performance statistics of all the methods can be found in Tables S2 and S4. Methods are referred to by their submission ID's which are provided in Table 1.

### **411 3.1 Analysis of macroscopic pK<sub>a</sub> predictions**



**Figure 2. RMSE and unmatched  $pK_a$  counts vs. submission ID plots for macroscopic  $pK_a$  predictions based on Hungarian matching.** Methods are indicated by submission IDs. RMSE is shown with error bars denoting 95% confidence intervals obtained by bootstrapping over challenge molecules. Submissions are colored by their method categories. Light blue colored database lookup methods are utilized as the null prediction method. QM methods category (navy) includes pure QM, QM+LEC, and QM+MM approaches. Lower bar plots show the number of unmatched experimental  $pK_a$ s (light grey, missing predictions) and the number of unmatched  $pK_a$  predictions (dark grey, extra predictions) for each method between pH 2 and 12. Submission IDs are summarized in Table 1. Submission IDs of the form  $nb\#\#\#$  refer to non-blinded reference methods computed after the blind challenge submission deadline. All others refer to blind, prospective predictions.

The performance of macroscopic  $pK_a$  predictions were analyzed by comparison to experimental  $pK_a$  values collected by the spectrophotometric method via numerical matching following the Hungarian method. Overall  $pK_a$  prediction performance was worse than we hoped. Fig. 2 shows RMSE calculated for each prediction method represented by their submission IDs. Other performance statistics are depicted in Fig. 3. In both figures, method categories were indicated by the color of the error bars. Statistics depicted in these figures can be found in Table S2. Prediction error ranged between 0.7 to 3.2  $pK_a$  units in terms of RMSE, while an RMSE between 2-3 log units was observed for the majority of methods (20 out of 38 methods). Only five methods achieved RMSE less than 1  $pK_a$  unit. One is QM method with COSMO-RS approach for solvation and linear empirical correction (*xvxzd* (DSD-BLYP-D3(B))/def2-TZVPD//PBEh-3c[DCOSMO-RS] + RRHO(GFN-xTB[GBSA]) + Gsolv(COSMO-RS[TZVPD]) and linear fit), and the remaining four are empirical prediction methods of LFER (*xmyhm* (ACD/pKa Classic), *nb007* (Schrödinger/EpiK Scan)) and QSPR/ML categories (*gyuhx* (Simulations Plus), *nb017* (MoKa)). These five methods with RMSE less than 1  $pK_a$  unit are also the methods that have the lowest MAE. *xmyhm* and *xvxzd* were the only two methods for which the upper 95% confidence interval of RMSE was lower than 1  $pK_a$  unit.

424 In terms of correlation statistics, many methods have good performance, although the ranking of methods changes according  
425 to  $R^2$  and Kendall's Tau. Therefore, many methods are indistinguishable from one another, considering the uncertainty of  
426 the correlation statistics. 32 out of 38 methods have  $R$  and Kendall's Tau higher than 0.7 and 0.6, respectively. 8 methods have  
427  $R^2$  higher than 0.9 and 6 methods have Kendall's Tau higher than 0.8. The overlap of these two sets are the following: *gyuhx* (Sim-

**Table 1. Submission IDs, names, category, and type for all the pKa prediction sets.** Reference calculations are labeled as *nb###*. The method name column lists the names provided by each participant in the submission file. The “type” column indicates if a submission was or a post-deadline reference calculation, denoted by “Blind” or “Reference” respectively. The methods in the table are grouped by method category and not ordered by performance.

Method Category	Method	Microscopic pKa (Type I) Submission ID	Macroscopic pKa (Type III) Submission ID	Submission Type	Ref.
DL	Substructure matches to experimental data in pKa OpenEye pKa Prospector Database v1.0	<i>5nm4j</i>	Null	[31]	
DL	OpenEye pKa-Prospector 1.0.0.3 with Analog Search ion identification algorithm	<i>pwn3m</i>	Null	[31]	
LFER	ACD/pKa GALAS (ACD/Percepta Kernel v1.6)	<i>v8qph</i>	<i>37xm8</i>	Blind	[32]
LFER	ACD/pKa Classic (ACD/Percepta Kernel, v1.6)		<i>xmyhm</i>	Blind	[33]
LFER	Epik Scan (Schrödinger v2017-4)		<i>nb007</i>	Reference	[25]
LFER	Epik Microscopic (Schrödinger v2017-4)	<i>nb008</i>	<i>nb010</i>	Reference	[25]
QSPR/ML	OpenEye Gaussian Process	<i>6tvf8</i>	<i>hyjnj</i>	Blind	[11]
QSPR/ML	OpenEye Gaussian Process Resampled		<i>q3pfj</i>	Blind	[11]
QSPR/ML	S+pKa (ADMET Predictor v8.5, Simulations Plus)	<i>hdijq</i>	<i>gyuhx</i>	Blind	[23]
QSPR/ML	Chemicalize v18.23 (ChemAxon MarvinSketch v18.23)		<i>nb015</i>	Reference	[34]
QSPR/ML	Moka v3.1.3	<i>nb016</i>	<i>nb017</i>	Reference	[21, 35]
QM	Adiabatic scheme with single point correction: SMD/M06-2X//6-311++G(d,p)//M06-2X/6-31+G(d) for bases and SMD/M06-2X//6-311++G(d,p)//M06-2X/6-31G(d) for acids + thermal corrections	<i>ko8yx</i>	<i>ryzue</i>	Blind	[36]
QM	Direct scheme with single point correction: SMD/M06-2X//6-311++G(d,p)//M06-2X/6-31G(d) for bases and SMD/M06-2X//6-311++G(d,p)//M06-2X/6-31G(d) for acids + thermal corrections	<i>w4z0e</i>	<i>xikp8</i>	Blind	[36]
QM	Adiabatic scheme: thermodynamic cycle that uses gas phase optimized structures for gas phase free energy and solution phase geometries for solvent phase free energy. SMD/M06-2X/6-31+G(d) for bases and SMD/M06-2X/6-31G(d) for acids + thermal corrections	<i>wcvnu</i>	<i>5byn6</i>	Blind	[36]
QM	Vertical scheme: thermodynamic cycle that uses only gas phase optimized structures to compute gas phase and solvation free energy. SMD/M06-2X/6-31+G(d) for bases and SMD/M06-2X/6-31G(d) for acids + Thermal corrections	<i>arcko</i>	<i>w4iyd</i>	Blind	[36]
QM	Direct scheme: solution phase free energy is determined by solution phase geometries without thermodynamic cycle SMD/M06-2X/6-31+G(d) for bases and SMD/M06-2X/6-31G(d) for acids + thermal corrections	<i>wexjs</i>	<i>y75vj</i>	Blind	[36]
QM + LEC	Jaguar (Schrödinger v2017-4)	<i>nb011</i>	<i>nb013</i>	Reference	[37]
QM + LEC	CPCM/B3LYP/6-311+G(d,p) and global fitting	<i>y4wws</i>	<i>35bdm</i>	Blind	[9]
QM + LEC	CPCM/B3LYP/6-311+G(d,p) and separate fitting for neutral to negative and for positive to neutral transformations	<i>qsicn</i>	<i>p0jba</i>	Blind	[9]
QM + LEC	EC-RISM/MP2/6-311+G(d,p)-P3NI-q-noThiols-2par	<i>kxzt</i>	<i>ds62k</i>	Blind	[38]
QM + LEC	EC-RISM/MP2/cc-pVTZ-P2-q-noThiols-2par	<i>ftc8w</i>	<i>2ii2g</i>	Blind	[38]
QM + LEC	EC-RISM/MP2/6-311+G(d,p)-P2-phi-all-2par	<i>ktpj5</i>	<i>nb001</i>	Blind*	[38]
QM + LEC	EC-RISM/MP2/6-311+G(d,p)-P2-phi-noThiols-2par	<i>wuuvc</i>	<i>nb002</i>	Blind*	[38]
QM + LEC	EC-RISM/MP2/6-311+G(d,p)-P3NI-phi-all-2par	<i>2umai</i>	<i>nb003</i>	Blind*	[38]
QM + LEC	EC-RISM/MP2/6-311+G(d,p)-P3NI-phi-noThiols-2par	<i>cm2yq</i>	<i>nb004</i>	Blind*	[38]
QM + LEC	EC-RISM/MP2/6-311+G(d,p)-P2-phi-all-1par	<i>z7fhp</i>	<i>nb005</i>	Blind*	[38]
QM + LEC	EC-RISM/MP2/6-311+G(d,p)-P3NI-phi-all-1par	<i>8toyp</i>	<i>nb006</i>	Blind*	[38]
QM + LEC	EC-RISM/MP2/cc-pVTZ-P2-phi-noThiols-2par	<i>epvmk</i>	<i>tjjd0</i>	Blind	[38]
QM + LEC	EC-RISM/MP2/cc-pVTZ-P2-phi-all-2par	<i>xnoe0</i>	<i>mkhqa</i>	Blind	[38]
QM + LEC	EC-RISM/MP2/cc-pVTZ-P3NI-phi-noThiols-2par	<i>4o0ia</i>	<i>mpwiy</i>	Blind	[38]
QM + LEC	EC-RISM/B3LYP/6-311+G(d,p)-P3NI-q-noThiols-2par	<i>nxaaw</i>	<i>ad5pu</i>	Blind	[38]
QM + LEC	EC-RISM/B3LYP/6-311+G(d,p)-P3NI-phi-noThiols-2par	<i>0xi4b</i>	<i>f0gew</i>	Blind	[38]
QM + LEC	EC-RISM/B3LYP/6-311+G(d,p)-P2-phi-noThiols-2par	<i>cwyk</i>	<i>np6b4</i>	Blind	[38]
QM + LEC	PCM/B3LYP/6-311+G(d,p)	<i>gdqeg</i>	<i>yc70m</i>	Blind	[38]
QM + LEC	COSMOtherm_FINE17 (COSMOtherm C30_1701, BP/TZVPD/FINE//BP/TZVP/COSMO)	<i>t8ewk</i>	<i>0hxtm</i>	Blind	[39, 40]
QM + LEC	DSD-BLYP-D3(BJ)/def2-TZVPD//PBEh-3c[DCOSMO-RS] + RRHO[GFN-xTB[GBSA]] + Gsolv(COSMO-RS[TZVPD]) and linear fit		<i>xvxzd</i>	Blind	[41]
QM + LEC	ReScosS conformations // DSD-BLYP-D3 reranking // COSMOtherm pKa: DSD-BLYP-D3(BJ)/def2-TZVPD// PBE-D3(BJ)/def2-TZVP/COSMO + RRHO[GFN-xTB + GBSA-water] + Gsolv[COSMO-RS(FINE17/TZVPD)] level and COSMOtherm pKa applied at the single conformer pair level (COSMOthermX17.0.5 release and BP-TZVPD-FINE-C30-1701 parameterization)	<i>eyetm</i>	<i>8xt50</i>	Blind	[41]
QM + LEC	ReScosS conformations // COSMOtherm pKa: DSD-BLYP-D3(BJ)/def2-TZVPD// PBE-D3(BJ)/def2-TZVP/COSMO + RRHO[GFN-xTB + GBSA-water] + Gsolv[COSMO-RS(FINE17/TZVPD)] level and COSMOtherm pKa was applied directly on the resulting conformer sets with at least 5% Boltzmann weights for each microspecies (COSMOthermX17.0.5 release and BP-TZVPD-FINE-C30-1701 parameterization)	<i>ccpmw</i>	<i>yqkga</i>	Blind	[41]
QM + MM	M06-2X/6-31G*(for bases) or 6-31+G*(for acids) for gas phase, solvation free energy using TI with explicit solvent and GAFF, solvation free energy of proton -265.6 kcal/mol	<i>0wfzo</i>		Blind	[42]
QM + MM	M06-2X/6-31G*(for bases) or 6-31+G*(for acids) for gas phase, solvation free energy using TI with explicit solvent and GAFF, solvation free energy of proton -271.88 kcal/mol	<i>z3btx</i>		Blind	
QM + MM	M06-2X/6-31G*(for bases) or 6-31+G*(for acids) + thermal state correction for gas phase, solvation free energy using TI with explicit solvent and GAFF, solvation free energy of proton -265.6 kcal/mol	<i>758j8</i>		Blind	
QM + MM	M06-2X/6-31G*(for bases) or 6-31+G*(for acids) + thermal state correction for gas phase, solvation free energy using TI with explicit solvent and GAFF, solvation free energy of proton -271.88 kcal/mol	<i>hgn83</i>		Blind	

\* Microscopic pKa submissions were blind, however, participant requested a correction after blind submission deadline for macroscopic pKa submissions. Therefore, these were assigned submission IDs in the form of *nb##*.

ulations Plus), *xvxzd* (DSD-BLYP-D3(BJ)/def2-TZVPD//PBEh-3c[DCOSMO-RS] + RRHO(GFN-xTB[GBSA]) + Gsolv(COSMO-RS[TZVPD]) and linear fit), *xmyhm* (ACD/pKa Classic), *ryzue* (Adiabatic scheme with single point correction: MD/M06-2X//6-311++G(d,p)//M06-2X/6-31+G(d) for bases and SMD/M06-2X//6-311++G(d,p)//M06-2X/6-31G(d) for acids + thermal corrections), and *5byn6* (Adiabatic scheme: thermodynamic cycle that uses gas phase optimized structures for gas phase free energy and solution phase geometries for solvent phase free energy. SMD/M06-2X/6-31+G(d) for bases and SMD/M06-2X/6-31G(d) for acids + thermal corrections). It is worth noting that the *ryzue* and *5byn6* are QM predictions without any empirical correction. Their high correlation and rank correlation coefficient scores signal that with an empirical correction their accuracy based performance could improve. Indeed, the participants have shown that this is the case in their own challenge analysis paper and achieved RMSE of 0.73 p*K<sub>a</sub>* units after the challenge [36].

Null prediction methods based on database lookup (*5nm4j* and *pwn3m*) had similar performance, with an RMSE of roughly 2.5 p*K<sub>a</sub>* units, an MAE of 1.5 p*K<sub>a</sub>* units, R<sup>2</sup> of 0.2 and Kendall's Tau of 0.3. Many methods were observed to have a prediction performance advantage over the null predictions shown in light blue in Fig. 2 and Fig. 3 considering all the performance metrics as a whole. In terms of correlation statistics, the null methods are the worst performers, except *0hxtm*. From the perspective of accuracy-based statistics (RMSE and MAE), only the top 10 methods were observed to have significantly lower errors than the null methods considering the uncertainty of error metrics expressed as 95% confidence intervals.

The distribution of macroscopic p*K<sub>a</sub>* prediction signed errors observed in each submission was plotted in Fig. 7A as ridge plots based on Hungarian matching. *2ii2g*, *f0gew*, *np64b*, *p0jba*, and *yc70m* tend to overestimate and *5byn6*, *ryzue*, and *w4iyd* tend to underestimate macroscopic p*K<sub>a</sub>* values.

There were four submissions of QM+LEC category that used the COSMO-RS implicit solvation model. While three of these achieved the lowest RMSE among QM-based methods (*xvxzd*, *yqkga*, and *8xt50*) [41], one of them showed the highest RMSE (*0hxtm* (COSMOtherm\_FINE17)) in SAMPL6 Challenge macroscopic p*K<sub>a</sub>* predictions. All four methods used COSMO-RS/FINE17 level to compute solvation free energies. The major difference between the three low-RMSE methods and the *0hxtm* seems to be the protocol for determining relevant conformations for each microstate. *xvxzd*, *yqkga*, and *8xt50* methods used semi-empirical tight binding (GFN-xTB) method and GBSA continuum solvation model for geometry optimization of conformers and followed up with high level single point energy calculations with solvation free energy (COSMO-RS(FINE17/TZVPD)) and rigid rotor harmonic oscillator (RRHO[GFN-xTB(GBSA)] corrections. *yqkga*, and *8xt50* methods selected conformations for each microstate with Relevant Solution Conformer Sampling and Selection (ReSCoSS) workflow. Conformations were clustered according to shape and lowest energy conformations from each cluster according to BP86/TZVP/COSMO single point energies in any of the 10 different COSMO-RS solvents were considered as relevant conformers. The ReSCoSS workflow was described more in detail by Pracht et al [41] *yqkga* method further filtered out conformers that have less than 5% Boltzmann weights at the DSD-BLYP-D3/def2-TZVPD + RRHO(GFNxTB) + COSMO-RS(fine) level. *xvxzd* method used MF-MD-GC//GFN-xTB workflow and used energy thresholds of 6 kcal/mol and 10 kcal/mol, for conformer and microstate selection. On the other hand, the conformational ensemble captured for each microstate seems to be more limited for *0hxtm* method, judging by the method description provided in the submission file (this participant did not publish an analysis of the results that they obtained for SAMPL6). *0hxtm* method reported that relevant conformations were computed with the COSMOconf 4.2 workflow which produced multiple relevant conformers for only the neutral states of SM18 and SM22. In contrast to *xvxzd*, *yqkga*, and *8xt50* methods, the *0hxtm* method also did not include a RRHO correction. Participants of the three low-RMSE methods report that capturing the chemical ensemble for each molecule including conformers and tautomers and high level QM calculations led to more successful macroscopic p*K<sub>a</sub>* prediction results and RRHO correction provided a minor improvement [41]. Comparing these results to other QM approaches in SAMPL Challenge also points to the advantage of the COSMO-RS solvation approach compared to other implicit solvent models.

In addition to the statistics related to the p*K<sub>a</sub>* value, we also analyzed missing or extra p*K<sub>a</sub>* predictions. Analysis of the p*K<sub>a</sub>* values with accuracy- and correlation-based error metrics was only possible after the assignment of predicted macroscopic p*K<sub>a</sub>*s to experimental p*K<sub>a</sub>*s through Hungarian matching, although this approach masks p*K<sub>a</sub>* prediction issues in the form of extra or missing macroscopic p*K<sub>a</sub>* predictions. To capture this form of prediction errors we reported the number of unmatched experimental p*K<sub>a</sub>*s (missing p*K<sub>a</sub>* predictions) and the number of unmatched predicted p*K<sub>a</sub>*s (extra p*K<sub>a</sub>* predictions) after Hungarian matching for each method. Both missing and extra p*K<sub>a</sub>* prediction counts were only considered for the pH range of 2-12 which was the limits of experimental assay. The lower subplot of Fig. 2 shows the total count of unmatched experimental or predicted p*K<sub>a</sub>*s for all the molecules in each prediction set. The order of submission IDs in the x-axis follows the RMSD based ranking so that the performance of each method from both p*K<sub>a</sub>* value accuracy and the number of p*K<sub>a</sub>*s can be viewed together. Presence of missing or extra macroscopic p*K<sub>a</sub>* predictions is a critical error because inaccuracy in predicting the correct number of macroscopic transitions shows that methods are failing to predict the correct set of charge states, i.e. failing to predict the correct

number of ionization states that can be observed between the specified pH range.

In challenge results, extra macroscopic  $pK_a$  predictions were found to be more common than missing  $pK_a$  predictions. In  $pK_a$  prediction evaluations, the accuracy of predicted ionization states within a pH range is usually neglected. When predictions are only evaluated for the accuracy of the  $pK_a$  value with numerical matching algorithms, larger number of predicted  $pK_a$ s lead to greater underestimation of prediction errors. Therefore, it is not surprising that methods are biased to predict extra  $pK_a$  values. The SAMPL6  $pK_a$  Challenge experimental data consists of 31 macroscopic  $pK_a$ s in total, measured for 24 molecules (6 molecules in the set have multiple  $pK_a$ s). Within the 10 methods with lowest RMSE, only the *xvxzd* method predicts too few  $pK_a$  values (2 unmatched out of 31 experimental  $pK_a$ s). All other methods that rank in the top 10 by RMSE have extra predicted  $pK_a$ s ranging from 1 to 13. Two prediction sets without any extra  $pK_a$  predictions and low RMSE are *8xt50* (ReSCoSS conformations // DSD-BLYP-D3 reranking // COSMOtherm pKa) and *nb015* (ChemAxon/Chemicalize).

### 3.1.1 Consistently well-performing methods for macroscopic $pK_a$ prediction

Methods ranked differently when ordered by different error metrics, although there were a couple of methods that consistently ranked at the top fraction. By using combinatorial criteria that take all multiple statistical metrics and unmatched  $pK_a$  counts into account, we identified a shortlist of consistently well-performing methods for macroscopic  $pK_a$  predictions, shown in Table 2. The criteria for selection were ranking in Top 10 according to RMSE, MAE,  $R^2$ , and Kendall's Tau and also having a combined unmatched  $pK_a$  (extra and missing  $pK_a$ s) count less than 8 (a third of the number of compounds). This resulted in a list of four methods that are consistently well-performing across all criteria.

Consistently well performing methods for macroscopic  $pK_a$  prediction included methods from all categories. Two methods of the QM+LEC category were *xvxzd* (DSD-BLYP-D3(BJ)/def2-TZVPD//PBEh-3c[DCOSMO-RS] + RRHO(GFN-xTB[GBSA]) + Gsolv(COSMO-RS[TZVPD]) and linear fit) and (*8xt50*) (ReSCoSS conformations // DSD-BLYP-D3 reranking // COSMOtherm pKa) and both used COSMO-RS approach. Empirical  $pK_a$  predictions with top performance were both proprietary software. From QSPR and LFER categories, *gyuhx* (Simulation Plus) and *xmyhbm* (ACD/pKa Classic) were consistently well-performing methods. The Simulation Plus  $pK_a$  prediction method consisted of 10 artificial neural network ensembles trained on 16,000 compounds for 10 classes of ionizable atoms. Atom type and the local molecular environment was how the ionization class of each atom was determined [43]. ACD/pKa Classic method which was trained on 17,000 compounds uses Hammett-type equations and tries to capture effects related to tautomeric equilibria, covalent hydration, resonance effects, and  $\alpha$ ,  $\beta$ -unsaturated systems [33].

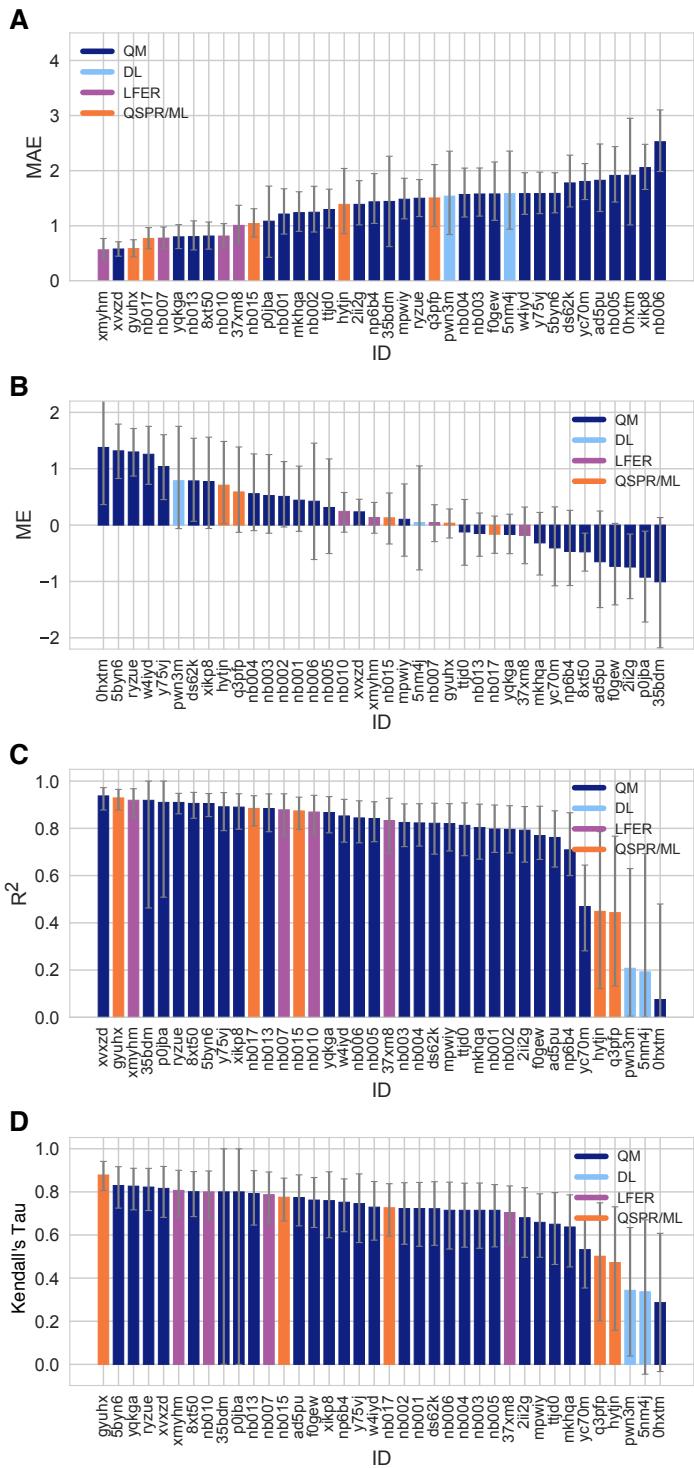
**Table 2. Four consistently well-performing prediction methods for macroscopic  $pK_a$  prediction based on consistent ranking within the Top 10 according to various statistical metrics.** Submissions were ranked according to RMSE, MAE,  $R^2$ , and  $\tau$ . Consistently well-performing methods were selected as the ones that rank in the Top 10 in each of these statistical metrics. These methods also have less than 2 unmatched experimental  $pK_a$ s and less than 7 unmatched predicted  $pK_a$ s according to Hungarian matching. Performance statistics are provided as mean and 95% confidence intervals.

Submission ID	Method Name	RMSE	MAE	$R^2$	Kendall's Tau ( $\tau$ )	Unmatched Exp. $pK_a$ Count	Unmatched Pred. $pK_a$ Count [2,12]
<i>xvxzd</i>	Full quantum chemical calculation of free energies and fit to experimental $pK_a$	0.68 [0.54, 0.81]	0.58 [0.45, 0.71]	0.94 [0.88, 0.97]	0.82 [0.68, 0.92]	2	4
<i>gyuhx</i>	S+pKa	0.73 [0.55, 0.91]	0.59 [0.44, 0.74]	0.93 [0.88, 0.96]	0.88 [0.8, 0.94]	0	7
<i>xmyhbm</i>	ACD/pKa Classic	0.79 [0.52, 1.03]	0.56 [0.38, 0.77]	0.92 [0.85, 0.97]	0.81 [0.68, 0.9]	0	3
<i>8xt50</i>	ReSCoSS conformations // DSD-BLYP-D3 reranking // COSMOtherm pKa	1.07 [0.78, 1.36]	0.81 [0.58, 1.07]	0.91 [0.84, 0.95]	0.80 [0.68, 0.89]	0	0

In Figure 4 prediction vs. experimental data correlation plots of macroscopic  $pK_a$  predictions with 4 consistently well-performing methods, a representative average method, and the null method(*5nm4j*). The representative method with average performance (*2ii2g* (EC-RISM/MP2/cc-pVTZ-P2-q-noThiols-2par)) was selected as the method with the highest RMSE below the median of all methods.

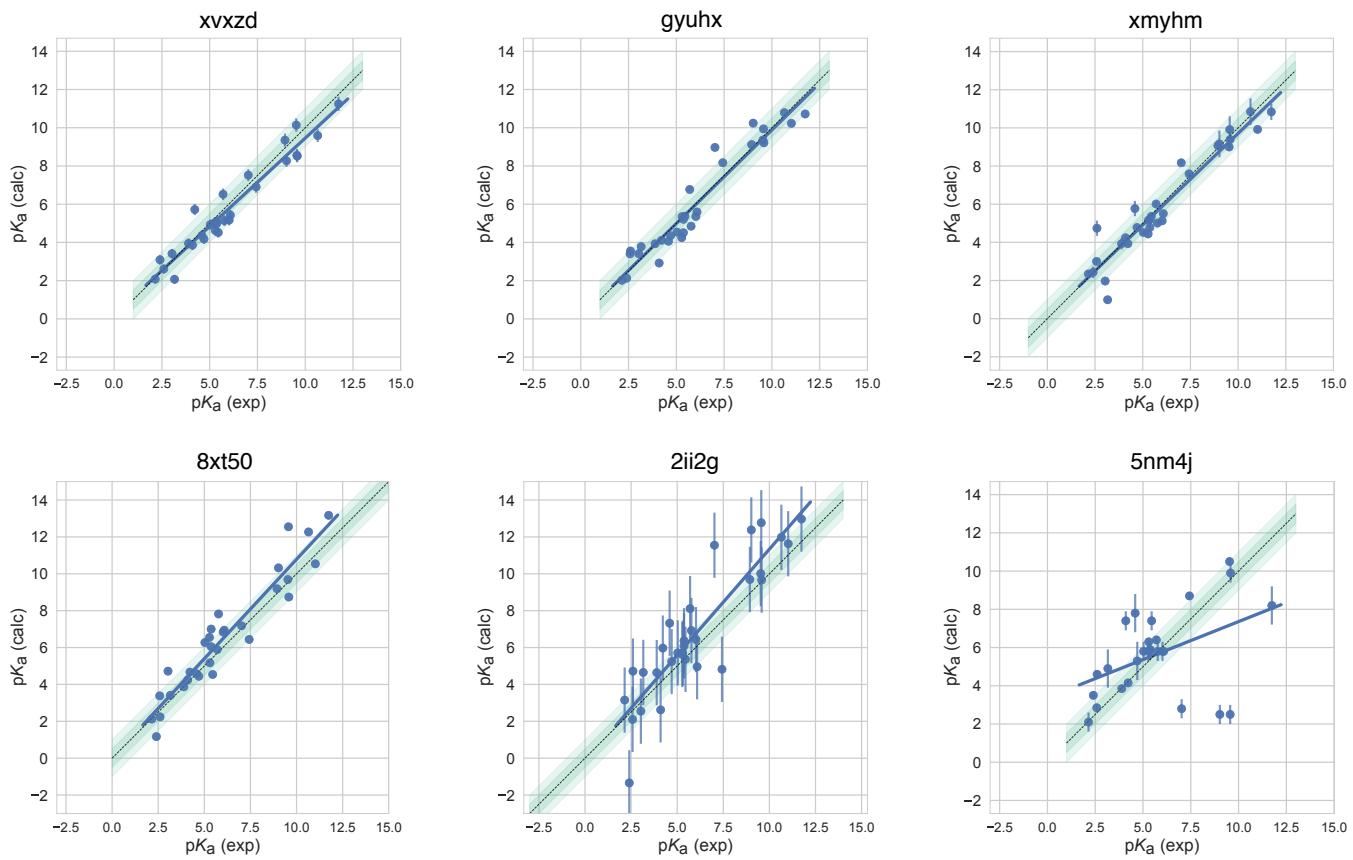
### 3.1.2 Which chemical properties are driving macroscopic $pK_a$ prediction failures?

In addition to comparing the performance of methods that participated in the SAMPL6 Challenge, we also wanted to analyze macroscopic  $pK_a$  predictions from the perspective of challenge molecules and determine whether particular compounds suffer from larger inaccuracy in  $pK_a$  predictions. The goal of this analysis is to provide insight on which molecular properties or moieties might be causing larger  $pK_a$  prediction errors. In Fig. 5, 2D depictions of the challenge molecules are presented with MAE calculated for their macroscopic  $pK_a$  predictions over all methods, based on Hungarian match. For multiprotic molecules,



**Figure 3. Additional performance statistics for macroscopic  $pK_a$  predictions based on Hungarian matching.** Methods are indicated by submission IDs. Mean absolute error (MAE), mean error (ME), Pearson's  $R^2$ , and Kendall's Rank Correlation Coefficient  $\tau$  are shown, with error bars denoting 95% confidence intervals obtained by bootstrapping over challenge molecules. Refer to Table 1 for submission IDs and method names. Submissions are colored by their method categories. Light blue colored database lookup methods are utilized as the null prediction method.

515 MAE was averaged over all the  $pK_a$ s. For the analysis of  $pK_a$  prediction accuracy observed for each molecule, MAE is a more  
 516 appropriate statistical value than RMSE for following global trends. This is because MAE value less sensitive to outliers than is



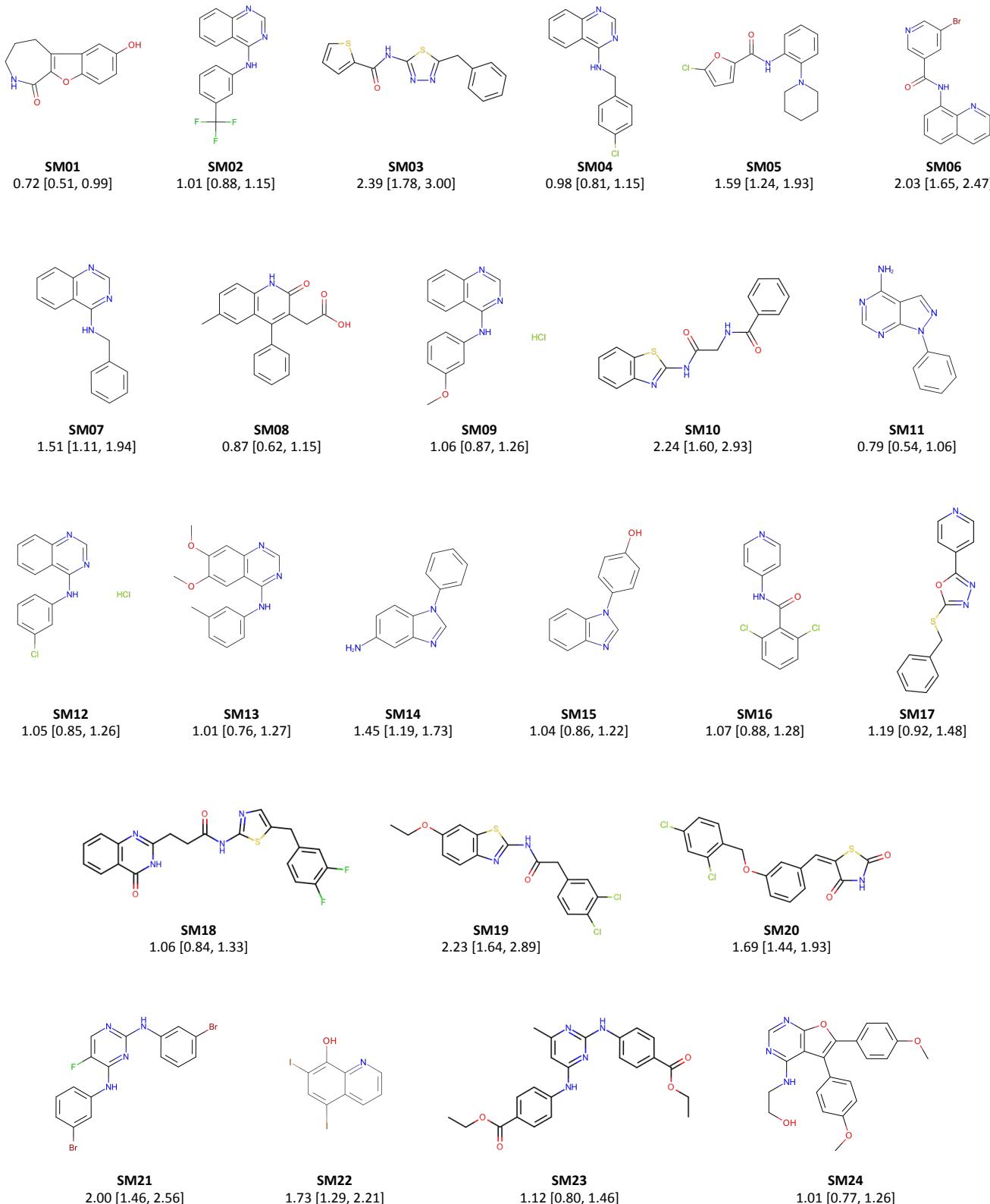
**Figure 4. Predicted vs. experimental value correlation plots of 4 consistently well-performing methods, a representative method with average performance (2ii2g), and the null method (5nm4j).** Dark and light green shaded areas indicate 0.5 and 1.0 units of error. Error bars indicate standard error of the mean of predicted and experimental values. Experimental  $pK_a$  SEM values are too small to be seen under the data points. EC-RISM/MP2/cc-pVTZ-P2-q-noThiols-2par method (2ii2g) was selected as the representative method with average performance because it is the method with the highest RMSE below the median.

#### RMSE.

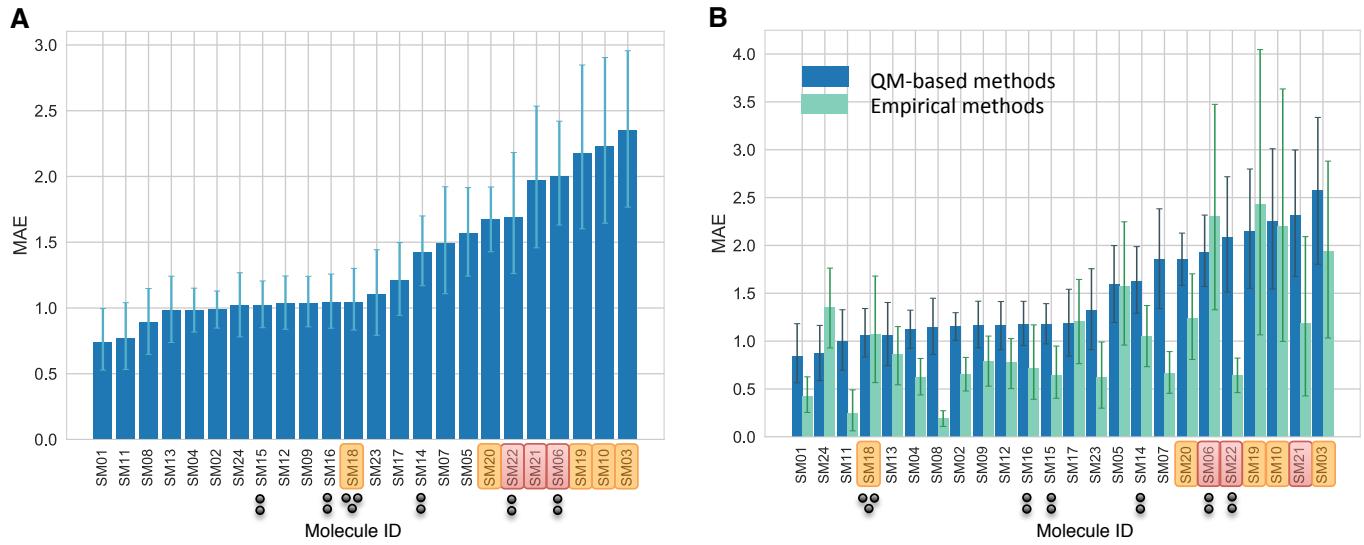
A comparison of the prediction accuracy of individual molecules is shown in Fig. 6. In Fig. 6A MAE each molecule is shown considering all blind predictions and reference calculations. A cluster of molecules marked orange and red have higher than average MAE. Molecules marked red (SM06, SM21, and SM22) are the only compounds in the SAMPL6 dataset with bromo or iodo groups and they suffered a macroscopic  $pK_a$  prediction error in the range of 1.7-2.0  $pK_a$  units in terms of MAE. Molecules marked orange (SM03, SM10, SM18, SM19, and SM20) have sulfur-containing heterocycles, and all these molecules except SM18 have MAE larger than 1.6  $pK_a$  unit. SM18 despite containing thiazole group has a low MAE. SM18 is the only compound with three experimental  $pK_a$ s and we suspect the presence of multiple experimental  $pK_a$ s could have a masking effect on the errors captured by MAE with Hungarian matching due to more pairing choices.

We analyzed MAE of each molecule for empirical (LFER and QSPR/ML) and QM-based physical methods (QM, QM+LEC, and QM+MM) separately for more insight. Fig. 6B shows that the difficulty of predicting  $pK_a$ s of the same subset of molecules was a trend conserved in the performance of physical methods. For QM-based methods too sulfur-containing heterocycles, amide next to aromatic heterocycles, compounds with iodo and bromo substitutions have lower  $pK_a$  prediction accuracy.

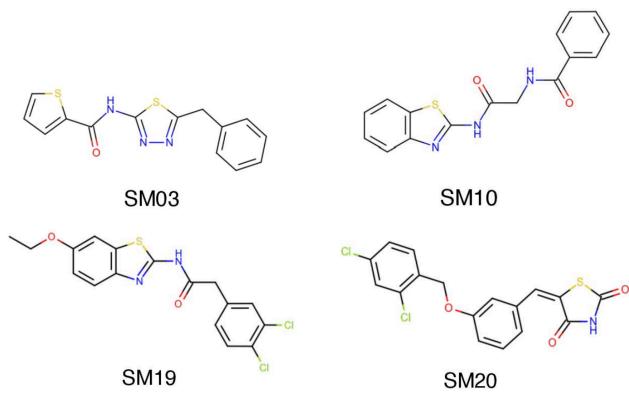
SAMPL6  $pK_a$  set consists of only 24 small molecules and lacks repeating instances of many moieties which limits our ability to determine with statistical significance which chemical substructures cause greater errors in  $pK_a$  predictions. Still, the trends observed in this challenge point to molecules with iodo, bromo, and sulfur-containing heterocycles with larger prediction errors of macroscopic  $pK_a$  value. We hope that reporting this observation will lead to the improvement of methods for similar compounds with such moieties.



**Figure 5. Molecules of SAMPL6 Challenge with MAE calculated for all macroscopic  $pK_a$  predictions.** MAE calculated considering all prediction methods indicate which molecules had the lowest prediction accuracy in SAMPL6 Challenge. MAE values calculated for each molecule include all the matched  $pK_a$  values. SM06, SM14, SM15, SM16, SM18, and SM22 were multiprotic. Hungarian matching algorithm was employed for pairing experimental and predicted  $pK_a$  values. MAE values are reported with 95% confidence intervals.

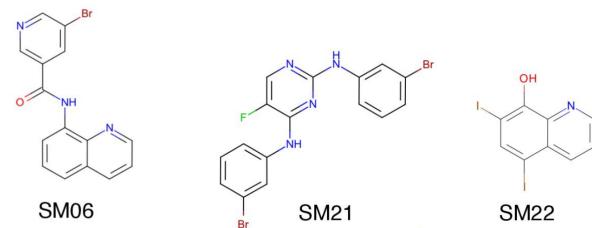


**C SAMPL6 molecules with sulfur-containing heterocycles**

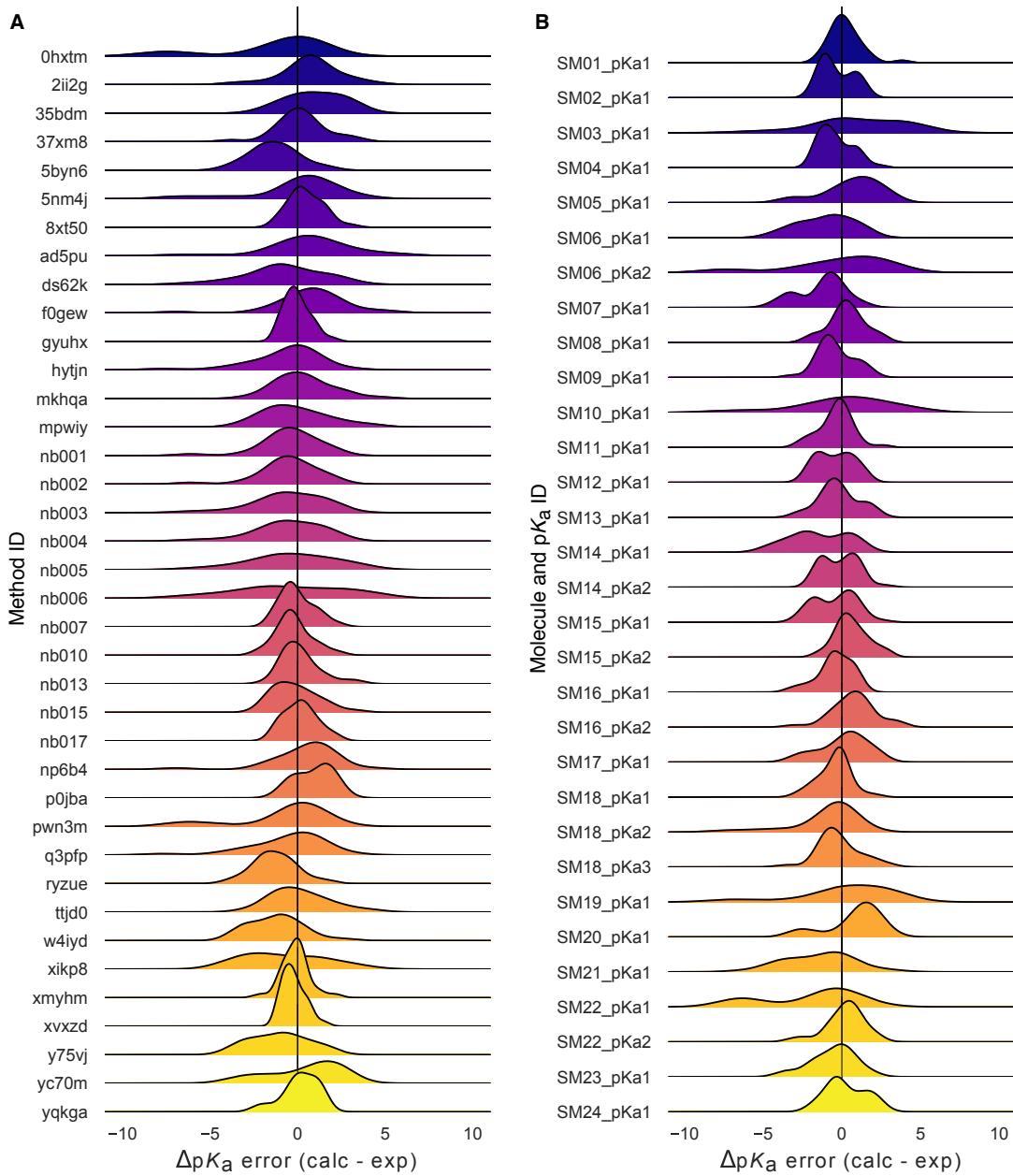


● 3 experimental  $pK_a$  values  
● 2 experimental  $pK_a$  values  
■ Sulfur-containing heterocycles  
■ Bromo and iodo groups

**D SAMPL6 molecules with bromo and iodo groups**



**Figure 6. Average prediction accuracy calculated over all prediction methods was lower for molecules with sulfur-containing heterocycles, bromo, and iodo groups.** (A) MAE calculated for each molecule as an average of all methods. (B) MAE of each molecule broken out by method category. QM-based methods (blue) include QM predictions with or without linear empirical correction. Empirical methods (green) include QSAR, ML, DL, and LFER approaches. (C) Depiction of SAMPL6 molecules with sulfur-containing heterocycles. (D) Depiction of SAMPL6 molecules with iodo and bromo groups.



**Figure 7. Macroscopic  $pK_a$  prediction error distribution plots show how prediction accuracy varies across methods and individual molecules.** (A)  $pK_a$  prediction error distribution for each submission for all molecules according to Hungarian matching. (B) Error distribution for each SAMPL6 molecule for all prediction methods according to Hungarian matching. For multiprotic molecules,  $pK_a$  ID numbers (pKa1, pKa2, and pKa3) were assigned in the direction of increasing experimental  $pK_a$  value.

535 We have also looked for correlation with molecular descriptors for finding other potential explanations for why macroscopic  
536  $pK_a$  prediction errors were larger for some molecules. While testing the correlation between errors and many molecular de-  
537 scriptors, it is important to keep the possibility of spurious correlations in mind. We haven't observed any significant correlation  
538 between numerical  $pK_a$  predictions and the descriptors we have tested. First, having more experimental  $pK_a$ s (Fig. 6A) did not  
539 seem to associate with worse  $pK_a$  prediction performance. But we need to keep in mind that there was a low representation  
540 of multiprotic compounds in the SAMPL6 set (5 molecules with 2 macroscopic  $pK_a$ s and one with 3 macroscopic  $pK_a$ ). Second,  
541 we checked the following other descriptors: amide group presence, molecular weight, heavy atom count, rotatable bond count,  
542 heteroatom count, heteroatom to carbon ratio, ring system count, maximum ring size, and the number of microstates (as enu-  
543 merated for the challenge). Correlation plots and  $R^2$  values can be seen in Fig. S2. We had suspected that  $pK_a$  prediction methods  
544 may be trained better for moderate values (4-10) than extreme values as molecules with extreme  $pK_a$ s are less likely to change  
545 ionization states close to physiological pH. To test this we look at the distribution of absolute errors calculated for all molecules  
546 and challenge predictions binned by experimental  $pK_a$  value 2  $pK_a$  unit increments. As can be seen in Fig. S3B, the value of true  
547 macroscopic  $pK_a$ s was not a factor affecting prediction error seen in SAMPL6 Challenge.

548 Fig. 7B is helpful to answer the question of "Are there molecules with consistently overestimated or underestimated  $pK_a$ s?".  
549 This ridge plots show the error distribution of each experimental  $pK_a$ . SM02\_pKa1, SM04\_pKa1, SM14\_pKa1, and SM21\_pKa1  
550 were underestimated, predicting lower protein affinity by more than 1  $pK_a$  unit by majority of the prediction methods. SM03\_pKa1,  
551 SM06\_pKa2, SM19\_pKa1, and SM20\_pKa1 were overestimated by the majority of the prediction methods by more than 1  $pK_a$  unit.  
552 SM03\_pKa1, SM06\_pKa2, SM10\_pKa1, SM19\_pKa1, and SM22\_pKa1 have the highest spread of errors and were less accurately  
553 predicted overall.

### 554 3.2 Analysis of microscopic $pK_a$ predictions using microstates determined by NMR for 8 molecules

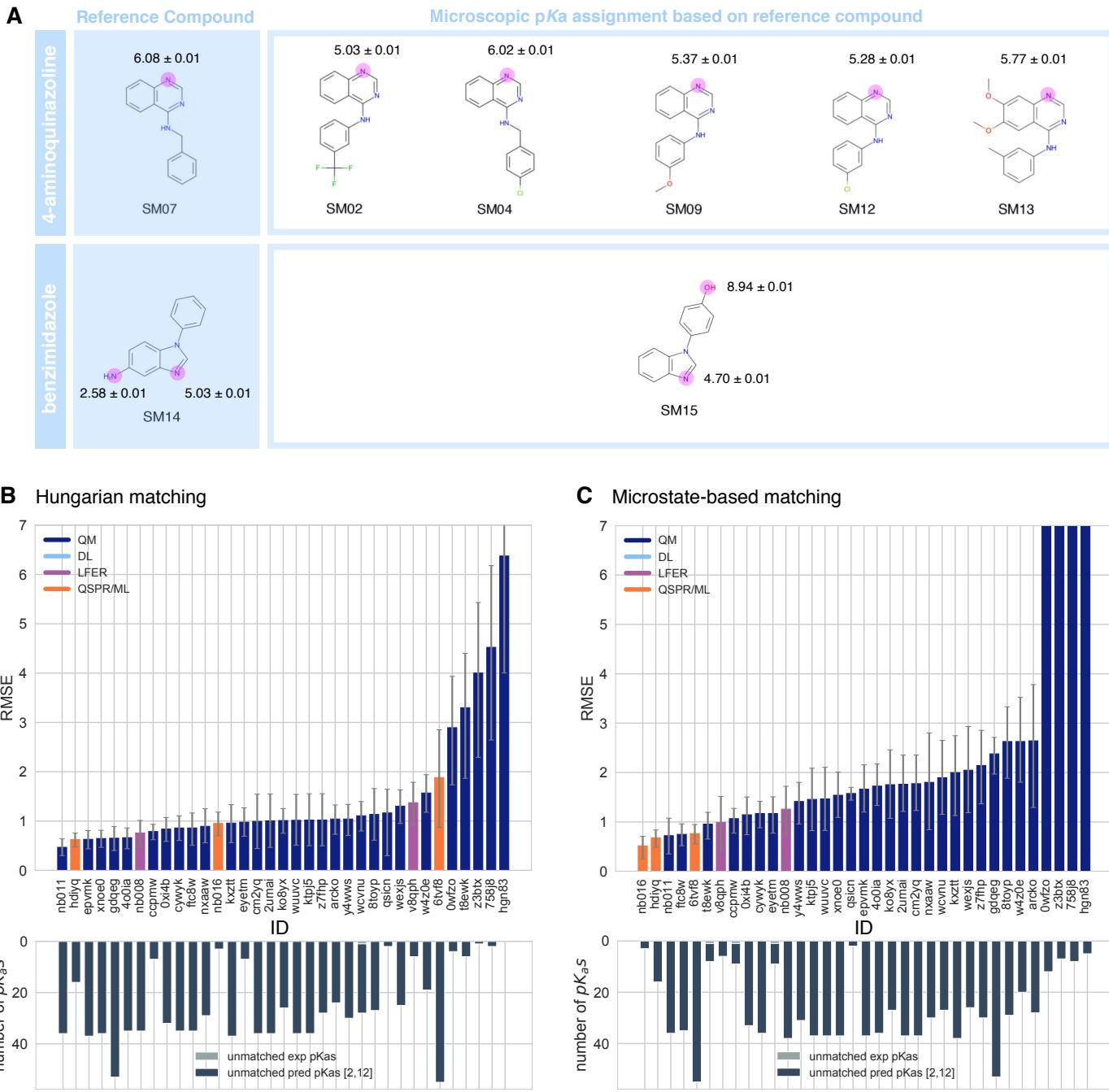
555 The common approach for analyzing microscopic  $pK_a$  prediction accuracy has been to compare it to experimental macroscopic  
556  $pK_a$  data, assuming experimental  $pK_a$ s describe titrations of distinguishable sides and, therefore, equal to microscopic  $pK_a$ s. But  
557 this typical approach fails to evaluate the methods at the microscopic level.

558 Analysis of microscopic  $pK_a$  predictions of the SAMPL6 Challenge was not straightforward due to the lack of experimental  
559 data with microscopic detail. For 24 molecules macroscopic  $pK_a$ s were determined with the spectrophotometric method. For 18  
560 molecules a single macroscopic titration was observed and for 6 molecules multiple experimental  $pK_a$ s were reported. For 18  
561 molecules with single experimental  $pK_a$ , it is probable that the molecules are monoprotic and therefore macroscopic  $pK_a$  value  
562 is equal to the microscopic  $pK_a$ . There is, however, no direct experimental evidence supporting this hypothesis but only the  
563 support from computational predictions. There is always the possibility that the macroscopic  $pK_a$  observed is the result of two  
564 different titrations overlapping closely with respect to pH if any charge state has more than one tautomer. We did not want to  
565 bias the blind challenge analysis with any prediction method. Therefore, we believe analyzing the microscopic  $pK_a$  predictions  
566 via Hungarian matching to experimental values with the assumption that the 18 molecules have a single titratable site is not the  
567 best approach. Instead, analysis at the level of macroscopic  $pK_a$ s is much more appropriate when a numerical matching scheme  
568 is the only option to evaluate predictions using macroscopic experimental data.

569 For a subset of the molecules in the dataset of 8 molecules, dominant microstates were inferred from NMR experiments.  
570 This dataset was extremely useful for guiding the assignment between experimental and predicted  $pK_a$  values based on mi-  
571 crostates. In this section, we present the performance evaluations of microscopic  $pK_a$  predictions for only the 8 compounds  
572 with experimentally determined dominant microstates.

#### 573 3.2.1 Microstate-based matching revealed errors masked by $pK_a$ value-based matching between experimental 574 and predicted $pK_a$ s

575 Comparing microscopic  $pK_a$  predictions directly to macroscopic experimental  $pK_a$  values with numerical matching can lead to  
576 underestimation of errors. To demonstrate how numerical matching often masks  $pK_a$  prediction errors we compared the per-  
577 formance analysis done by Hungarian matching to that from microstate-based matching for 8 molecules presented in Fig. 8A. RMSE  
578 calculated for microscopic  $pK_a$  predictions matched to experimental values via Hungarian matching is shown in Fig. 8B, while  
579 Fig. 8C shows RMSE calculated via microstate-based matching. The Hungarian matching incorrectly leads to significantly lower  
580 RMSE compared to microstate-based matching. The reason is that the Hungarian matching assigns experimental  $pK_a$  values to  
581 predicted  $pK_a$  values only based on the closeness of the numerical values, without consideration of the relative population of  
582 microstates and microstate identities. Because of that, a microscopic  $pK_a$  value that describes a transition between very low  
583 population microstates (high energy tautomers) can be assigned to the experimental  $pK_a$  if it has the closest  $pK_a$  value. This is



**Figure 8. NMR determination of dominant microstates allowed in-depth evaluation of microscopic pKa predictions of 8 compounds.**

**A** Dominant microstate sequence of two compounds (SM07 and SM14) were determined by NMR [7]. Based on these reference compounds dominant microstates of 6 other derivative compounds were inferred and experimental pKa values were assigned to titratable groups with the assumption that only the dominant microstates have significant contributions to the experimentally observed pKa. **B** RMSE vs. submission ID and unmatched pKa vs. submission ID plots for the evaluation of microscopic pKa predictions of 8 molecules by Hungarian matching to experimental macroscopic pKas. **C** RMSE vs. submission ID and unmatched pKa vs. submission ID plots showing the evaluation of microscopic pKa predictions of 8 molecules by microstate-based matching between predicted microscopic pKas and experimental macroscopic pKa values. Submissions *0wfzo*, *z3btx*, *758j8*, and *hgn83* have RMSE values bigger than 10 pKa units which are beyond the y-axis limits of subplot **C** and **B**. RMSE is shown with error bars denoting 95% confidence intervals obtained by bootstrapping over challenge molecules. Lower bar plots show the number of unmatched experimental pKas (light grey, missing predictions) and the number of unmatched pKa predictions (dark grey, extra predictions) for each method between pH 2 and 12. Submission IDs are summarized in Table 1.

584 not helpful because, in reality, the microscopic  $pK_a$ s that influence the observable macroscopic  $pK_a$  the most are the ones with  
585 higher populations (transitions between low energy tautomers).

586 The number of unmatched predicted microscopic  $pK_a$ s is shown in the lower bar plots of Fig. 8B and C, to emphasize the large  
587 number of microscopic  $pK_a$  predictions submitted by many methods. In the case of microscopic  $pK_a$ , the number of unmatched  
588 predictions does not indicate an error in the form of an extra predicted  $pK_a$ , because the spectrophotometric experiments do  
589 not capture all microscopic  $pK_a$ s theoretically possible (transitions between all pairs of microstates that are 1 proton apart).  
590  $pK_a$ s of transitions to and from very high energy tautomers are very hard to measure by experimental methods, including the  
591 most sensitive methods like NMR. Prediction of extra microscopic  $pK_a$ s can cause underestimation of prediction errors when  
592 numerical matching algorithms such as Hungarian matching are used. We also checked how often Hungarian matching led to  
593 the correct matches between predicted and experimental  $pK_a$  in terms of the microstate pairs, i.e. how often the microstate  
594 pair of the Hungarian match recapitulates the dominant microstate pair of the experiment. The overall accuracy of microstate  
595 pair matching was found to be low for SAMPL6 Challenge submission. Fig. S4 shows that for most methods the predicted  
596 microstate pair selected by the Hungarian match did not correspond to the experimentally determined microstate pair. This  
597 means the lower RMSE results obtained from Hungarian matching are low for the wrong reason. This problem could be avoided  
598 by matching experimental and predicted values on the basis of microstate IDs.

599 Unfortunately, we are only able to perform this more reliable microstate-based matching analysis for a subset of com-  
600 pounds. The conclusions in this section are only about eight compounds with limited diversity. This subset is composed of six  
601 4-aminoquinazoline molecules and two with benzimidazole scaffolds, for a total of 10  $pK_a$  values. The sequences of dominant  
602 microstates for SM07 and SM14 were determined by NMR experiments directly [7] and dominant microstates of their derivatives  
603 were inferred by taking them as reference (Fig. 8). Although we believe that microstate-based evaluation is more informative, the  
604 lack of a large experimental dataset limits the conclusions to a very narrow chemical diversity. Still, microstate-based matching  
605 revealed errors masked by  $pK_a$  value-based matching between experimental and predicted  $pK_a$ s.

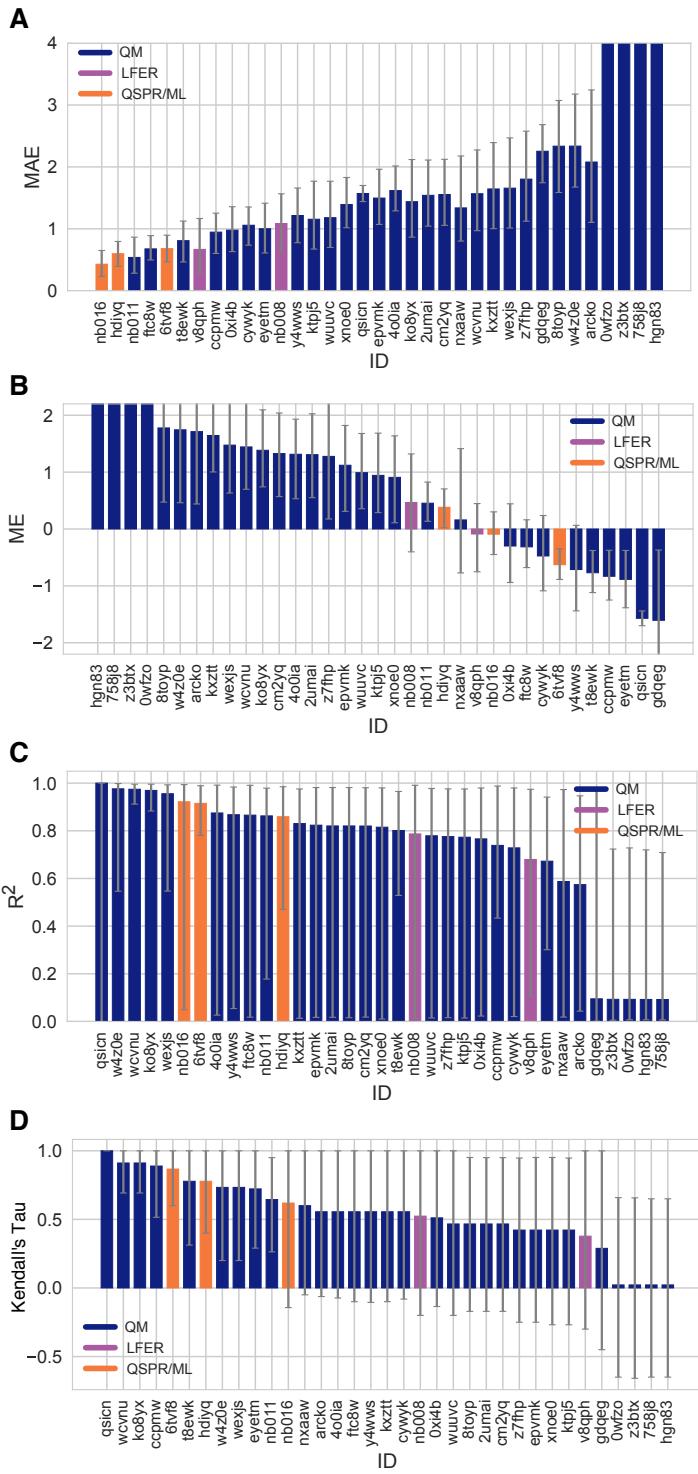
### 606 3.2.2 Accuracy of $pK_a$ predictions evaluated by microstate-based matching

607 Both accuracy and correlation based statistics were calculated for predicted microscopic  $pK_a$  values after microstate-based  
608 matching. RMSE, MAE, ME,  $R^2$ , and Kendall's Tau results of each method are shown in Fig. 8C and Fig. 9. A table of the calculated  
609 statistics can be found in Table S4. Due to small number of data points in this set, correlation-based statistics have large uncer-  
610 tainties and thus have less utility for distinguishing better performing methods. Therefore, we focused more on accuracy based  
611 metrics for the analysis of microscopic  $pK_a$ s than correlation based metrics. In terms of accuracy of microscopic  $pK_a$  value, all  
612 three QSPR/ML based methods (*nb016* (MoKa), *hdiyq* (Simulations Plus), *6tvf8* (OE Gaussian Process)), three QM-based methods  
613 (*nb011* (Jaguar), *ftc8w* (EC-RISM/MP2/cc-pVTZ-P2-q-noThiols-2par), *t8ewk* (COSMOlogic\_FINE17)), and one LFER method (*v8qph*  
614 (ACD/pKa GALAS)) achieved RMSE lower than 1  $pK_a$  unit. The same 6 methods also have the lowest MAE.

### 615 3.2.3 Evaluation of dominant microstate prediction accuracy

616 For many computational chemistry approaches including structure-based modeling of protein-ligand interactions, predicting  
617 the ionization state and the exact position of protons is needed to establish what to include in the modeled system. This is  
618 why in addition to being able to predict  $pK_a$  values accurately, we need  $pK_a$  prediction methods to be able to capture micro-  
619 scopic protonation states accurately. Even when the predicted  $pK_a$  value is very accurate, the predicted protonation site can be  
620 wrong. Therefore, we assessed if methods participating in the SAMPL6  $pK_a$  Challenge were predicting correctly the sequence of  
621 dominant microstates, i.e. dominant tautomers of each charge state observed between pH 2 and 12.

622 Fig. 10 shows how well methods perform for predicting the dominant microstate, as analyzed for eight compounds with  
623 required experimental data. The dominant microstate sequence is essentially the sequence of states that are most visible  
624 experimentally, due to their higher fractional population and relative free energy within the tautomers at each charge. To extract  
625 the dominant tautomers predicted for the sequence of ionization states of each method, the relative free energy of microstates  
626 were first calculated at reference pH 0 [27]. Then to determine the dominant microstate at each formal charge, we have selected  
627 the lowest energy tautomer for each ionization state based on the relative microstate free energies calculated at pH 0. The  
628 choice of reference pH is not important, as relative free energy difference between tautomers of the same charge is always  
629 constant with respect to pH. This analysis was done only for the charges -1, 0, 1, and 2 which was the charge range captured  
630 by NMR experiments. Then predicted and experimental dominant microstates were compared for each charge to calculate the  
631 fraction of correctly predicted dominant tautomers. This value is reported as the dominant microstate accuracy for all charges  
632 (Fig. 10A). Many of the methods which participated the challenge made errors in predicting the dominant microstate. 10 QM



**Figure 9. Additional performance statistics for microscopic  $pK_a$  predictions for 8 molecules with experimentally determined dominant microstates.** Microstate-based matching was performed between experimental  $pK_a$  values and predicted microscopic  $pK_a$  values. Mean absolute error (MAE), mean error (ME), Pearson's  $R^2$ , and Kendall's Rank Correlation Coefficient Tau ( $\tau$ ) are shown, with error bars denoting 95% confidence intervals obtained by bootstrapping over challenge molecules. Methods are indicated by their submission IDs. Submissions are colored by their method categories. Refer to Table 1 for submission IDs and method names. Submissions 0wfzo, z3btx, 758j8, and hgn83 have MAE and ME values bigger than 10  $pK_a$  units which are beyond the y-axis limits of subplots A and B. A large number and wide variety of methods have a statistically indistinguishable performance based on correlation statistics (C and D), in part because of the relatively small dynamic range the small size of the set of 8 molecules.

and 3 QSPR/ML methods did not make any mistakes in dominant microstate predictions, although, they are expected to make mistakes in the relative population of tautomers (free energy difference between microstates) as reflected by the  $pK_a$  value errors. While all the participating QSPR/ML methods showed good performance in dominant microstate prediction, LFER and some QM methods made mistakes. The accuracy of the predicted dominant neutral tautomers was perfect for all methods, except *qsicn* (Fig. 10B). But errors in predicting the major tautomer of charge +1 were much more frequent. 22 out of 35 prediction sets made at least one error in prediction the lowest energy tautomer with +1 charge. We didn't include ionization states with charges -1 and +2 in this assessment because we had only one compound with these charges in the dataset. Nevertheless, errors in predicting the dominant tautomers seem to be a bigger problem for charged tautomers than the neutral tautomer.

Only eight compounds had data on the sequence of dominant microstates. Therefore conclusions on the performance of methods in terms of dominant tautomer prediction are limited to this limited chemical diversity (benzimidazole and 4-aminoquinazoline derivatives). We present this analysis as a prototype of how microscopic  $pK_a$  predictions should be evaluated. Hopefully, future evaluation can be done with larger experimental datasets following the strategy we demonstrated here in order to reach broad conclusions about which methods are better for capturing dominant microstates and ratios of tautomers. Even if experimental microscopic  $pK_a$  measurement data is not available, experimental dominant tautomer determinations are still informative for assessing computational predictions.

The most frequent misprediction was the major tautomer of the SM14 cationic form, as shown in Fig. 10. This figure shows the accuracy of the predicted dominant microstate calculated for individual molecules and for charge states 0 and +1, averaged over all prediction methods. SM14, the molecule that exhibits the most frequent error in the predicted dominant microstate, has two experimental  $pK_a$  values that were 2.4  $pK_a$  units apart, and we suspect that could be a contributor to the difficulty of predicting microstates accurately. Other molecules are monoprotic (4-aminoquinazolines) or their experimental  $pK_a$  values are very well separated (SM14, 4.2  $pK_a$  units). It would be very interesting to expand this assessment to a larger variety of drug-like molecules to discover for which structures tautomer predictions are more accurate and for which structure computational predictions are not as reliable.

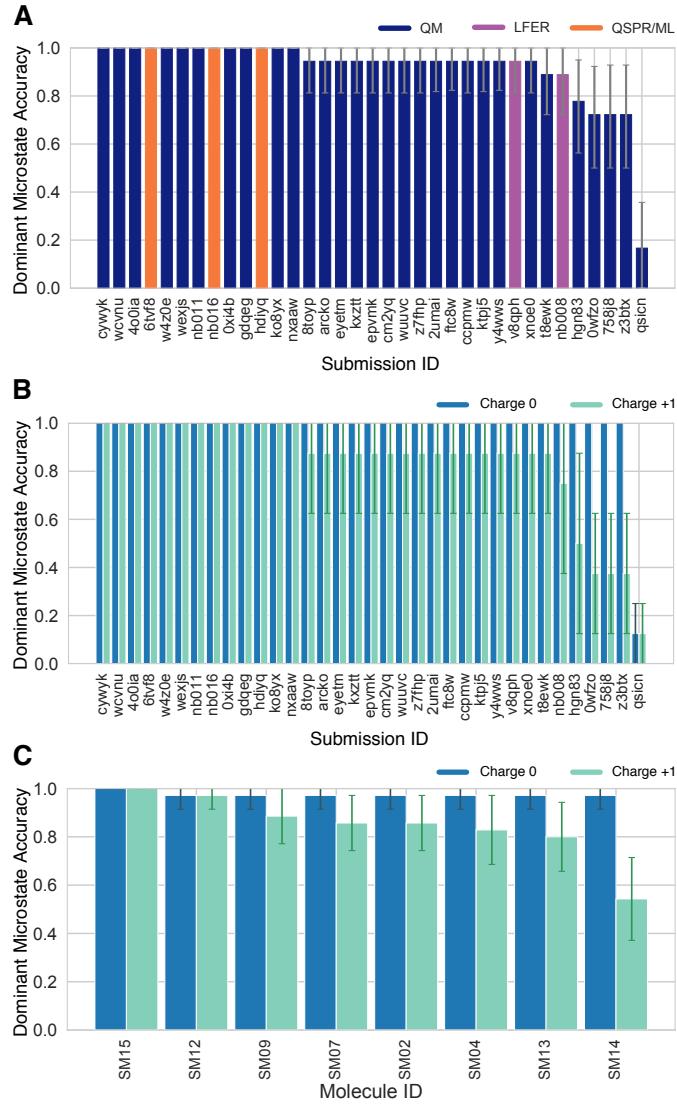
### 3.2.4 Consistently well-performing methods for microscopic $pK_a$ predictions

We have determined different criteria for determine consistently top-performing predictions of microscopic  $pK_a$  than macroscopic  $pK_a$ : having perfect dominant microstate prediction accuracy, unmatched  $pK_a$  count of 0, and ranking in the top 10 according to RMSE and MAE. Correlation statistics were not found to have utility for discriminating performance due to large uncertainties in these statistics for a small dataset of 10  $pK_a$  values. Unmatched predicted  $pK_a$  count was also not a consideration, since experimental data was only informative for the  $pK_a$  between dominant microstates and did not capture the all possible theoretical transitions between microstate pairs. Table 3 reports six methods that have consistent good performance according to many metrics, although evaluated only for the 8 molecule set due to limitations of the experimental dataset. Six methods were divided evenly between methods of QSPR/ML category and QM category. *nb016* (MoKa), *hdlyq* (Simulations Plus), and *6tvf8* (OE Gaussian Process) were QSPR and ML methods that performed well. *nb011* (Jaguar), *0xi4b*(EC-RISM/B3LYP/6-311+G(d,p)-P2-phi-noThiols-2par), and *cywyk* (EC-RISM/B3LYP/6-311+G(d,p)-P2-phi-noThiols-2par) were QM predictions with linear empirical corrections with good performance with microscopic  $pK_a$  predictions.

The Simulations Plus  $pK_a$  prediction method is the only method that appeared to be consistently well performing in both the assessment for macroscopic and microscopic  $pK_a$  prediction (*gyuhx* and *hdlyq*). However it is worth noting that two methods that were consistently in top-performing methods list for macroscopic  $pK_a$  predictions lacked equivalent submissions of their underlying microscopic  $pK_a$  predictions and therefore could not be evaluated at the microstate level. These methods were (ACD/Classic  $pK_a$ ) and *xvxzd*(DSD-BLYP-D3(BJ)/def2-TZVPD//PBEh-3c[DCOSMO-RS] + RRHO(GFN-xTB[GBSA]) + Gsolv(COSMO-RS[TZVPD]) and linear fit).

## 3.3 How do $pK_a$ prediction errors impact protein-ligand binding affinity predictions?

$pK_a$  predictions provide a key input for computational modeling of protein-ligand binding with physical methods. The SAMPL6  $pK_a$  Challenge focused focused only on small molecule  $pK_a$  prediction and  $pK_a$  prediction accuracy observed can effect modeling of ligands. Many affinity prediction methods such as docking, MM/PBSA, MM/GBSA, absolute or alchemical relative free energy calculation methods predict the affinity of the ligand to a receptor in a fixed protonation state. These models strictly depend on  $pK_a$  predictions for determining possible protonation states of the ligand in the aqueous environment and in a protein complex, as well as the free energy penalty to reach those states [3]. The accuracy of  $pK_a$  predictions can become a limitation for the performance of physical models that try to capture molecular association.



**Figure 10. Some methods predicted the sequence of dominant tautomers inaccurately.** Prediction accuracy of the dominant microstate of each charged state was calculated using the dominant microstate sequence determined by NMR for 8 molecules as reference. **(A)** Dominant microstate accuracy vs. submission ID plot was calculated considering all the dominant microstates seen in the experimental microstate dataset of 8 molecules. **(B)** Dominant microstate accuracy vs. submission ID plot was generated considering only the dominant microstates of charge 0 and +1 seen in the 8 molecule dataset. The accuracy of each molecule is broken out by the total charge of the microstate. **(C)** Dominant microstate prediction accuracy calculated for each molecule averaged over all methods. In **(B)** and **(C)**, the accuracy of predicting the dominant neutral tautomer is showed in blue and the accuracy of predicting the dominant +1 charged tautomer is shown in green. Error bars denoting 95% confidence intervals obtained by bootstrapping.

In terms of the ligand protonation states, there are two ways in which  $pK_a$  prediction errors can influence the prediction accuracy for protein-ligand binding free energies as depicted in Fig. 11. The first scenario is when a ligand is present in aqueous solution in multiple protonation states (Fig. 11A). When only the minor aqueous protonation state contributes to protein-ligand complex formation, the overall binding free energy ( $\Delta G_{bind}$ ) needs to be calculated as the sum of binding free energy of the minor state and the protonation penalty of that state ( $\Delta G_{prot}$ ).  $\Delta G_{prot}$  is a function of pH and  $pK_a$ . A 1 unit of error in  $pK_a$  value would lead to 1.36 kcal/mol error in overall binding free energy if the protonation state with the minor population binds the protein. The equations in Fig. 11A show the calculation of overall free energy.

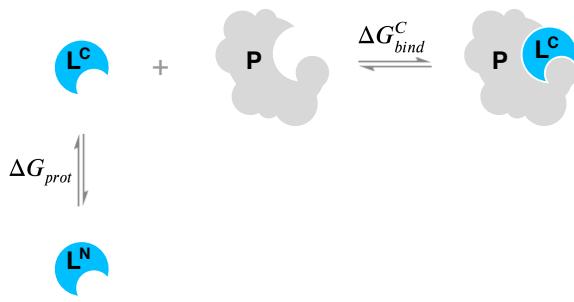
In addition to multiple protonation states being present in the aqueous environment, multiple charge states can contribute to complex formation (Fig. 11B). Then, the overall free energy of binding needs to include a Multiple Protonation States Correction (MPSC) term ( $\Delta G_{corr}$ ). MPSC is a function of pH, aqueous  $pK_a$  of the ligand, and the difference between the binding free energy

**Table 3. Top-performing methods for microscopic pK<sub>a</sub> predictions based on consistent ranking within the Top 10 according to various statistical metrics calculated for 8 molecule dataset.** Performance statistics are provided as mean and 95% confidence intervals. Submissions that rank in the Top 10 according to RMSE and MAE, and have perfect dominant microstate prediction accuracy were selected as consistently well-performing methods. Correlation-based statistics ( $R^2$ , and Kendall's Tau), although reported in the table, were excluded from the statistics used for determining top-performing methods. This was because correlation-based statistics were not very discriminating due to narrow dynamic range and the small number of data points in the 8 molecule dataset with NMR-determined dominant microstates.

Submission ID	Method Name	Dominant Microstate Accuracy	RMSE	MAE	R <sup>2</sup>	Kendall's Tau	Unmatched Exp. pK <sub>a</sub> Count	Unmatched Pred. pK <sub>a</sub> Count [2,12]
nb016	MoKa	1.0 [1.0, 1.0]	0.52 [0.25, 0.71]	0.43 [0.23, 0.65]	0.92 [0.05, 0.99]	0.62 [-0.14, 1.00]	0	3
hd1yq	S+pKa	1.0 [1.0, 1.0]	0.68 [0.49, 0.83]	0.60 [0.39, 0.80]	0.86 [0.47, 0.98]	0.78 [0.40, 1.00]	0	16
nb011	Jaguar	1.0 [1.0, 1.0]	0.72 [0.35, 1.07]	0.54 [0.28, 0.86]	0.86 [0.18, 0.98]	0.64 [0.26, 0.95]	0	36
6tvf8	OE Gaussian Process	1.0 [1.0, 1.0]	0.76 [0.55, 0.95]	0.68 [0.46, 0.90]	0.92 [0.78, 0.99]	0.87 [0.6, 1.00]	0	55
0xi4b	EC-RISM/B3LYP/6-311+G(d,p)-P3NI-phi-noThiols-2par	1.0 [1.0, 1.0]	1.15 [0.75, 1.50]	0.98 [0.63, 1.36]	0.77 [0.02, 0.98]	0.51 [-0.14, 1.00]	0	33
cywyk	EC-RISM/B3LYP/6-311+G(d,p)-P2-phi-noThiols-2par	1.0 [1.0, 1.0]	1.17 [0.88, 1.41]	1.06 [0.74, 1.35]	0.73 [0.02, 0.98]	0.56 [-0.08, 1.00]	0	36

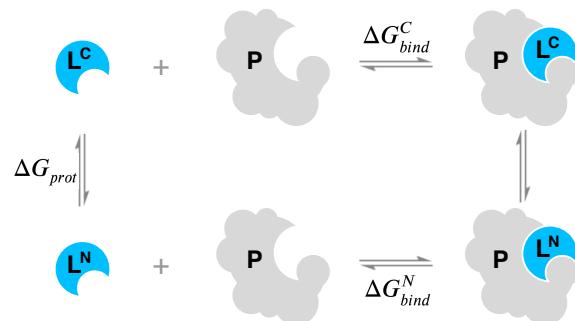
692 of charged and neutral species ( $\Delta G_{bind}^C - \Delta G_{bind}^N$ ) as shown in Fig. 11B.

### A When only the minor protonation state can bind to the protein



$$\Delta G_{bind} = \Delta G_{bind}^C + RT(pH - pK_a) \ln(10)$$

### B When multiple protonation states can bind to the protein



**Figure 11. Aqueous pK<sub>a</sub> of the ligand can influence overall protein-ligand binding affinity.** **A** When only the minor aqueous protonation state contributes to protein-ligand complex formation, overall binding free energy ( $\Delta G_{bind}$ ) needs to be calculated as the sum of binding affinity of the minor state and the protonation penalty of that state. **B** When multiple charge states contribute to complex formation, the overall free energy of binding includes a multiple protonation states correction (MPSC) term ( $\Delta G_{corr}$ ). MPSC is a function of pH, aqueous pK<sub>a</sub> of the ligand, and the difference between the binding free energy of charged and neutral species ( $\Delta G_{bind}^C - \Delta G_{bind}^N$ ).

693 Using the equations in Fig. 11B we can model the true MPSC ( $\Delta G_{corr}$ ) value with respect to the difference between pH and  
 694 the pK<sub>a</sub> of the ligand, to see when this value has significant impact to the overall binding free energy. In Fig. 12, the true MPSC  
 695 value that needs to be added to  $\Delta G_{bind}^N$  is shown for ligands with varying binding affinity difference between protonation states  
 696 ( $\Delta G = \Delta G_{bind}^C - \Delta G_{bind}^N$ ) and varying free energy of binding difference between the protonation states. Fig. 12A shows the case  
 697 of a monoprotic base in which the charged state has a lower affinity than the neutral state. Solid lines show the true correction.  
 698 In cases where the pK<sub>a</sub> is lower than the pH, the correction factor disappears as the ligand fully populates the neutral state  
 699 ( $\Delta G_{bind} = \Delta G_{bind}^N$ ). As the pK<sub>a</sub> value gets larger than the pH, the charged state is populated more and  $\Delta G_{corr}$  increases to approach  
 700  $\Delta G$ . It is interesting to note that the pH-pK<sub>a</sub> range over which  $\Delta G_{corr}$  changes. It is often assumed that for a basic ligand if  
 701 the pK<sub>a</sub> of a ligand is more than 2 units higher than the pH, then only 1% of the population is in the neutral state and it is safe  
 702 to approximate the overall binding affinity with  $\Delta G_{bind}^C$ . Based on the magnitude of the relative free energy difference between  
 703 ligand protonation states, this assumption is not always correct. As seen in Fig. 12A, responsive region of  $\Delta G_{corr}$  can span 3 pH

units for a system with  $\Delta\Delta G = 1 \text{ kcal/mol}$  or 5 pH units for a system with  $\Delta\Delta G = 4 \text{ kcal/mol}$ . This highlights that the range of  $pK_a$  values that impact binding affinity predictions is wider than previously appreciated. Molecules with  $pK_a$ s several units away from the physiological pH can still impact the overall binding affinity significantly due to the MPSC.

Despite the need to capture the contributions of multiple protonations states by including MPSC in binding affinity calculations, inaccurate  $pK_a$  predictions can lead to errors in  $\Delta G_{corr}$  and overall free energy of binding prediction. In Fig. 12A dashed lines show predicted  $\Delta G_{corr}$  based on  $pK_a$  error of -1 units. We have chosen a  $pK_a$  error of 1 unit as this is the average performance expected from the  $pK_a$  prediction methods based on the SAMPL6 Challenge. Underestimated  $pK_a$  causes underestimated  $\Delta G_{corr}$  and overestimated affinities (i.e. too negative binding free energy) for a varying range of pH -  $pK_a$  values depending on binding affinity difference between protonation states( $\Delta\Delta G$ ). In Fig. 12B dashed lines show how the magnitude of the absolute error caused by calculating  $\Delta G_{corr}$  with an inaccurate  $pK_a$  varies with respect to pH. Different colored lines show simulated results with varying binding free energy differences between protonation states. For a system whose charged state has higher binding free energy than the neutral state ( $\Delta\Delta G = 2 \text{ kcal/mol}$ ), the absolute error caused by underestimated  $pK_a$  by 1 unit only can be up to 0.9 kcal/mol. For a system whose charged state has even lower affinity (more positive binding free energy) than the neutral state ( $\Delta\Delta G = 4 \text{ kcal/mol}$ ), the absolute error caused by underestimated  $pK_a$  by 1 unit only can be up to 1.2 kcal/mol. The magnitude of errors contributing to overall binding affinity is too large to be neglected. Improving the accuracy of small molecule  $pK_a$  prediction methods can help to minimize the error in predicted MPSC.

With the current level of  $pK_a$  prediction accuracy as observed in SAMPL6 Challenge, is it advantageous to include MPSC in affinity predictions that may include errors caused by  $pK_a$  predictions? We provide a comparison of the two choices to answer this question: (1) Neglecting MPSC completely and assuming overall binding affinity is captured by  $\Delta G_{bind}^N$ , (2) including MPSC with potential error in overall affinity calculation. The magnitude of error caused by Choice 1 (ignoring MPSC) is depicted as a solid line in Fig. 12B and the magnitude of error caused by MPSC computed with inaccurate  $pK_a$  is depicted as dashed lines. What is the best strategy? Error due to choice 1 is always larger than error due to choice 2 for all pH- $pK_a$  values. In this scenario including MPSC improves overall binding affinity prediction. The error caused by the inaccurate  $pK_a$  is smaller than the error caused by neglecting MPSC.

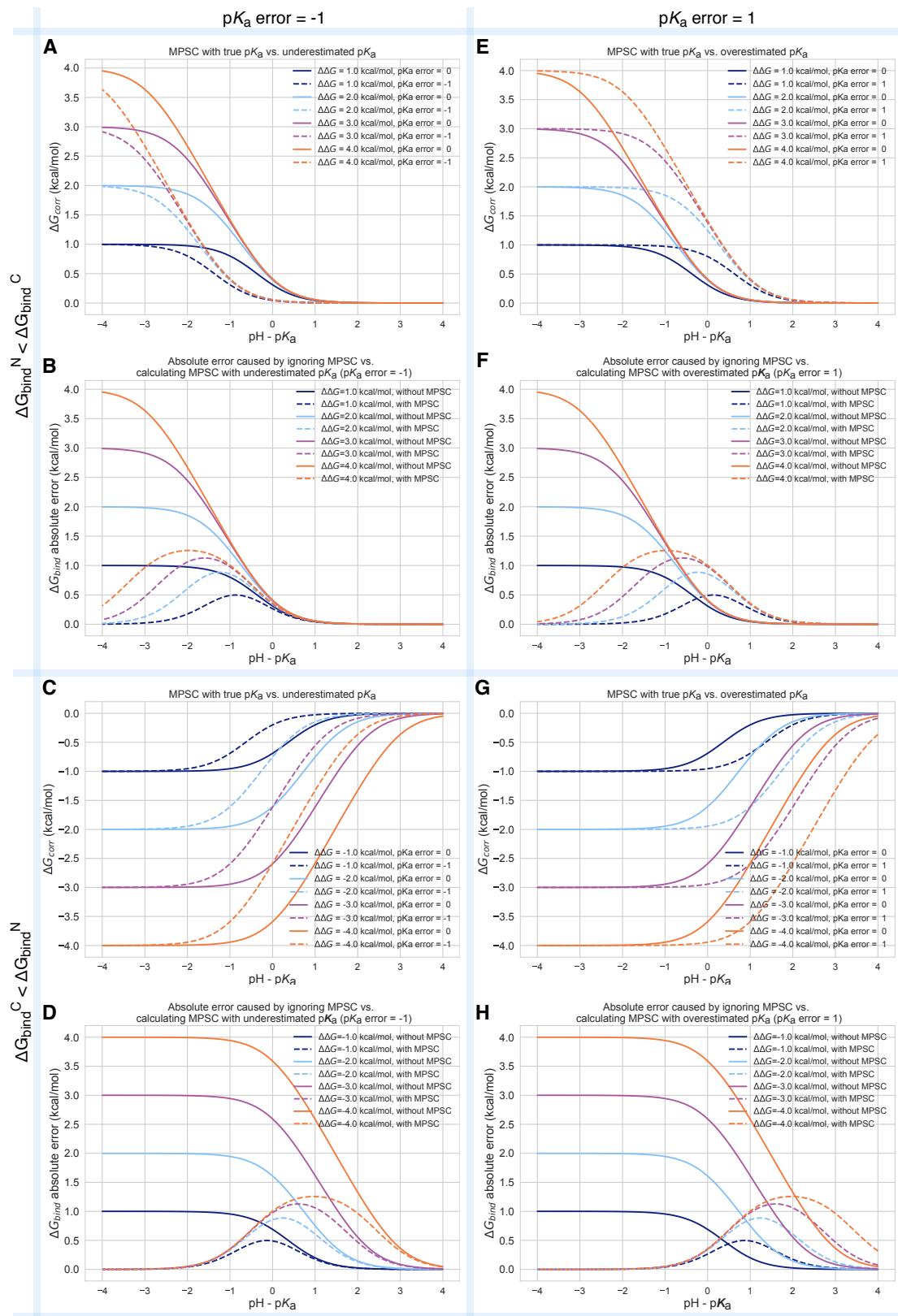
We can also ask ourselves whether or not an MPSC calculated based on an inaccurate  $pK_a$  should be included in binding affinity predictions in different circumstances such as underestimated or overestimated  $pK_a$  values and charged states with higher or lower affinities than the neutral states. We tried to capture these 4 circumstances in four quadrants of Fig. 12. In the case of overestimated  $pK_a$  values (Fig. 12E-H) it can be seen that for the most of the pH- $pK_a$  range it is more advantageous to include the predicted MPSC in affinity calculations, except a smaller window where the opposite choice would be more advantageous. For instance, for the system with  $\Delta\Delta G = 2 \text{ kcal/mol}$  and overestimated  $pK_a$  (Fig. 12E) for the pH- $pK_a$  region between -0.5 and 2, including predicted  $\Delta G_{corr}$  causes more error than ignoring MPSC.

In practice, we normally do not know the exact magnitude or the direction of the error of our predicted  $pK_a$ . Therefore, using simulated MPSC error plots to decide when to include MPSC in binding affinity predictions is not possible. However, based on the analysis of a case with 1 unit of  $pK_a$  error, including MPSC correction would be more often than not helpful in improving binding affinity predictions. The detrimental effect of  $pK_a$  inaccuracy is still significant. Hopefully, future improvements in  $pK_a$  prediction methods can improve the accuracy of MPSC and binding affinity predictions of ligands which have multiple protonation states that contribute to aqueous or complex populations. Being able to predict  $pK_a$ s with 0.5 units accuracy would significantly help the binding affinity models to incorporate more accurate MPSC terms.

### 3.4 Take-away lessons from SAMPL6 $pK_a$ Challenge

The SAMPL6  $pK_a$  Challenge showed that in general  $pK_a$  prediction performance of computational methods is lower than expected for drug-like molecules. Our expectation prior to the blind challenge was that well-developed methods would achieve prediction errors as low as 0.5  $pK_a$  units and reliable predictions of charge and tautomer states. There are many factors that complicate predicting  $pK_a$  values of drug-like molecules: multiple titratable sites, tautomerization, frequent presence of heterocycles, and extended conjugation patterns, as well as a high number of rotatable bonds, and the possibility of intramolecular hydrogen bonds. Macroscopic  $pK_a$  predictions have not yet reached experimental accuracy (Inter-method variability of macroscopic  $pK_a$  measurements can be around 0.5  $pK_a$  units [22]). There was not a single method in the SAMPL6 Challenge that achieved RMSE around 0.5 or lower for macroscopic  $pK_a$  predictions for the 24 molecule set of kinase inhibitor fragment-like molecules. Lower RMSE values were observed in the microscopic  $pK_a$  evaluation section of this study for some methods; however, the 8 molecule set used for that analysis poses a very limited dataset to reach conclusions about general expectations for drug-like molecules.

As the majority of experimental data was in the form of macroscopic  $pK_a$  values, we had to adopt a numerical matching



**Figure 12. Inaccuracy of  $pK_a$  prediction ( $\pm 1$  unit) affects the the accuracy of MPSC and overall protein-ligand binding free energy calculation in varying amounts based on aqueous  $pK_a$  value and relative binding affinity of individual protonation states ( $\Delta\Delta G = \Delta G_{bind}^C - \Delta G_{bind}^N$ ). All calculations are made for 25°C, and a ligand with a single basic titratable group. **A, C, E, and G** show MPSC ( $\Delta G_{corr}$ ) calculated with true vs. inaccurate  $pK_a$ . **B, D, F, and H** show the comparison of the absolute error to  $\Delta G_{bind}$  caused by ignoring the MPSC completely (solid lines) vs. calculating MPSC based in inaccurate  $pK_a$  value (dashed lines). These plots provide guidance on when it is beneficial to include MPSC correction based on  $pK_a$  error, pH -  $pK_a$ , and  $\Delta\Delta G$ .**

754 algorithm (Hungarian matching) to pair predicted and experimental values to calculate performance statistics of macroscopic  
755  $pK_a$  predictions. Accuracy, correlation, and extra/missing  $pK_a$  prediction counts were the main metrics for macroscopic  $pK_a$   
756 evaluations. An RMSE range of 0.7 to 3.2  $pK_a$  units was observed. Only five methods achieved RMSE between 0.7-1  $pK_a$  units,  
757 while an RMSE between 1.5-3 log units was observed for the majority of methods. All four methods of the LFER category and three  
758 out of 5 QSPR/ML methods achieved RMSE less than 1.5  $pK_a$  units. All the QM methods that achieved this level of performance  
759 included linear empirical corrections to rescale and unbias their  $pK_a$  predictions.

760 Based on the consideration of multiple error metrics, we compiled a shortlist of consistently-well performing methods for  
761 macroscopic  $pK_a$  evaluations. Two methods from QM+LEC methods, one QSPR/ML, two empirical methods achieved consistent  
762 performance according to many metrics. The common features of the two empirical methods were their large training sets  
763 (16000-17000 compounds) and being commercial prediction models.

764 There were four submissions of QM-based methods that utilized COSMO-RS implicit solvation model. While three of these  
765 achieved the lowest RMSE among QM-based methods (*xvxzd*, *yqkga*, and *8xt50*) [41], one of them showed the highest RMSE  
766 (*0hxtm* (COSMOtherm\_FINE17)). The comparison of these methods indicates that capturing conformational ensemble of mi-  
767 crostates, high level QM calculations, and RRHO corrections contribute to better macroscopic  $pK_a$  predictions. Linear empirical  
768 corrections applied QM calculations improved results, especially when the linear correction is calibrated for an experimental  
769 dataset using the same level of theory as the deprotonation free energy predictions (as in *xvxzd*). This challenge also points to  
770 the advantage of COSMO-RS solvation approach compared to other implicit solvent models.

771 Molecules that posed greatest difficulty for  $pK_a$  predictions were determined by comparing the macroscopic  $pK_a$  prediction  
772 accuracy of each molecules averaged over all methods submitted to the challenge.  $pK_a$  prediction errors were higher for com-  
773 pounds with sulfur-containing heterocycles, iodo, and bromo groups. This trend was also conserved when only QM-based  
774 methods were analyzed. SAMPL6  $pK_a$  dataset consisted of only 24 small molecules which limited our ability to statistically con-  
775 firm this conclusion, however, we believe it is worth reporting molecular features that coincided with larger errors even if we  
776 can not evaluate the reason for these failures.

777 Utilizing a numerical matching algorithm to pair experimental and predicted macroscopic  $pK_a$  values was a necessity, how-  
778 ever, this approach did not capture all aspects of prediction errors. Computing the number of missing or extra  $pK_a$  predictions  
779 remaining after Hungarian matching, provided a window of observing macroscopic  $pK_a$  prediction errors such as the number  
780 of macroscopic transitions or ionization states expected in a pH interval. In  $pK_a$  evaluation studies it is very typical to just focus  
781 on  $pK_a$  value errors evaluated after matching, and to ignore  $pK_a$  prediction errors that the matching protocol can not capture.  
782 Ignored prediction errors include predicting missing or extra  $pK_a$ s and failing to predict the correct charge states. SAMPL6  $pK_a$   
783 Challenge results showed sporadic presence of missing  $pK_a$  predictions and very frequent case of extra  $pK_a$  predictions. Both  
784 indicate failures to capture the correct ionization states. The traditional way of evaluating  $pK_a$ s that only focuses on the  $pK_a$   
785 value error after some sort of numerical match between predictions and experimental values may have motivated these types  
786 of errors as there would be no penalty for missing a macroscopic deprotonation and predicting an extra one. This problem does  
787 not seem to be specific to any method category.

788 We used the eight molecule subset of SAMPL6 compounds with NMR-based dominant microstate sequence information to  
789 demonstrate the advantage of evaluating  $pK_a$  prediction on the level of microstates. Comparison of statistics computed for the  
790 8 molecule dataset by Hungarian matching and microstate-based matching showed how Hungarian matching, despite being the  
791 optimal matching algorithm, can mask errors in  $pK_a$  predictions. Errors computed by microstate-based matching were larger  
792 compared to numerical matching algorithms in terms of RMSE. Microscopic  $pK_a$  analysis with numerical matching algorithms  
793 may mask errors due to the higher number of guesses made. Numerical matching based on  $pK_a$  values also ignores information  
794 regarding the relative population of states. Therefore, it can lead to  $pK_a$ s defined between very low energy microstate pairs to  
795 be matched to the experimentally observable  $pK_a$  between microstates of higher populations. Of course, the predicted  $pK_a$   
796 value could be correct however the predicted microstates would be wrong. Such mistakes caused by Hungarian matching were  
797 observed frequently in SAMPL6 results and therefore we decided microstate-based matching of  $pK_a$ values provides a more  
798 realistic picture of method performance.

799 Some QM and LFER methods made mistakes in predicting the dominant tautomers of the ionization states. Dominant tau-  
800 tomer prediction seemed to be a more prominent problem for charged tautomers than the neutral tautomer. The easiest way  
801 to extract dominant microstate sequence from predictions is to calculate the relative free energy of microstates at any reference  
802 pH, and determining the lowest energy state in each ionization state. Errors in dominant microstate predictions were very rare  
803 for neutral tautomers but more frequent in cationic tautomers with +1 charge of the 8 molecule set. SM14 was the molecule  
804 with the lowest dominant microstate prediction accuracy, while dominant microstates predictions for SM15 were perfect for all

805 molecules. SM14 and SM15 both have two experimental  $pK_a$ s and benzimidazole scaffold. The difference between them is the  
806 distance between the experimental  $pK_a$  values which is smaller for SM14. These results make sense from the perspective of  
807 relative free energies of microstates. Closer  $pK_a$  values mean that the free energy difference between different microstates are  
808 smaller for SM14, and therefore any error in predicting the relative free energy of tautomers is more likely to cause reordering of  
809 relative populations of microstates and impact the accuracy of dominant microstate predictions. It would have been extremely  
810 informative to evaluate the tautomeric ratios and relative free energy predictions of microstates, however, experimental data  
811 was missing for this approach.

812 The overall assessment of SAMPL6  $pK_a$  Challenge captured non-stellar performance for microscopic and macroscopic  $pK_a$   
813 predictions which can be detrimental to the accuracy of protein-ligand affinity predictions and other pH-dependent physico-  
814 chemical property predictions such as distribution coefficients, membrane permeability, and solubility. Protein-ligand binding  
815 affinity predictions rely on  $pK_a$  predictions in two ways: determination of relevant aqueous microstates and the free energy  
816 penalty to reach these states. Microscopic  $pK_a$  predictions with better accuracy are needed for accurate incorporation of mul-  
817 tiple protonation state correction (MPSC) to overall binding affinity calculations. We simulated the effect of overestimating or  
818 underestimating  $pK_a$  of a ligand by one unit on overall binding affinity prediction for a ligand where both cation and neutral  
819 states contribute to binding affinity.  $pK_a$  prediction error of this magnitude (assuming dominant tautomers were predicted cor-  
820 rectly) could cause up to 0.9 and 1.2 kcal/mol error in overall binding affinity when the binding affinity of protonation states  
821 are 2 or 4 kcal/mol different, respectively. For the case of 4 kcal/mol binding affinity difference between protonation states the  
822 pH- $pK_a$  range that the error would be larger than 0.5 kcal/mol surprisingly spans around 3.5 pH units. We demonstrated that the  
823 range of pH- $pK_a$  value that MPSC needs to be incorporated in binding affinity predictions can be wider than the widely assumed  
824 range of 2 pH units, based on the affinity difference between protonation states. At the level of 1 unit  $pK_a$  error incorporating  
825 MPSC would improve binding affinity predictions more often than not. If microscopic  $pK_a$  could be predicted with 0.5  $pK_a$  units  
826 of accuracy, MPSC calculations would be much more reliable.

827 There are multiple factors to consider when deciding which  $pK_a$  prediction method to utilize. These factors include the  
828 accuracy of microscopic and macroscopic  $pK_a$  values, the accuracy of the number and the identity of ionization states predicted  
829 within the experimental pH interval, the accuracy of microstates predicted within the experimental pH interval, the accuracy of  
830 tautomeric ratio (i.e. relative free energy between microstates), how costly is the calculation in terms of time and resources, and  
831 whether one has access to software licenses that might be required.

832 All of the top-performing empirical methods were developed as commercial software that require licenses to run, and there  
833 were not any open-source alternatives for empirical  $pK_a$  predictions. Since then two publications reported open-source machine  
834 learning-based  $pK_a$  prediction methods, however, one can only predict the most acidic or most basic macroscopic  $pK_a$  values  
835 of a molecule [44] and the second one is only trained for predicting  $pK_a$  values of monoprotic molecules [45]. Recently a  $pK_a$   
836 prediction methodology was published that describes a mixed approach of semi-empirical QM calculations and machine learn-  
837 ing that can predict macroscopic  $pK_a$ s of both mono-and polyprotic species [46]. The authors reported RMSE of 0.85 for the  
838 retrospective analysis performed on the SAMPL6 dataset.

### 839 **3.5 Suggestions for future blind challenge design and evaluation of $pK_a$ predictions**

840 This analysis helped us understand the current state of the field and led to many lessons. We believe the highest benefit can be  
841 achieved if further iteration so of small molecule  $pK_a$  prediction challenges can be organized, creating motivation for improving  
842 protonation state prediction methods for drug-like molecules. In future challenges, it is desirable to increase chemical diversity  
843 to cover more of common scaffolds [47] and functional groups [48] seen in drug-like molecules, and gradually increasing the  
844 complexity of molecules.

845 Future challenges should promote stringent evaluation for  $pK_a$  prediction methods from the perspective of microscopic  $pK_a$   
846 and microstate predictions. It is necessary to assess the capability of  $pK_a$  prediction methods to capture the free energy profile of  
847 microstates of multiprotic molecules. This is critical because  $pK_a$  predictions are often utilized to determine relevant protonation  
848 states and tautomers of small molecules that must be captured in other physical modeling approaches, such as protein-ligand  
849 binding affinity or distribution coefficient predictions. Different tautomers can have different binding affinities and partition  
850 coefficients.

851 In this paper, we demonstrated how experimental microstate information can guide the analysis further than the typical  $pK_a$   
852 evaluation approach that has been used so far. The traditional  $pK_a$  evaluation approach focuses solely on the numerical error  
853 of the  $pK_a$  values and neglects the difference between macroscopic and microscopic  $pK_a$  definitions. This is mainly caused by  
854 the lack of  $pK_a$  datasets with microscopic detail. To improve  $pK_a$  and protonation state predictions of multiprotic molecules it

855 is necessary to embrace the difference between macroscopic and microscopic  $pK_a$  definitions and select strategies for experi-  
856 mental data collection and prediction evaluation accordingly. In SAMPL6 Challenge the analysis was limited by the availability of  
857 experimental microscopic data as well. As usual macroscopic  $pK_a$  values were abundant (24 molecules) and limited data on mi-  
858 croscopic states was available (8 molecules), although the later opened new avenues for evaluation. For future blind challenges  
859 for multiprotic compounds, striving to collect experimental datasets with microscopic  $pK_a$ s would be very beneficial. Benchmark  
860 datasets of microscopic  $pK_a$ s are currently missing. This limits the improvement of  $pK_a$  and tautomer prediction methods for  
861 multiprotic molecules. If the collection of experimental microscopic  $pK_a$ s is not possible due to time and resource cost of such  
862 NMR experiments, at least supplementing the more automated macroscopic  $pK_a$  measurements with NMR-based determination  
863 of the dominant microstate sequence or tautomeric ratios of each ionization state can create very useful benchmark datasets.  
864 This supplementary information can allow microstate-based assignment between experimental and predicted  $pK_a$ s and a more  
865 realistic assessment of method performance.

866 If the only available experimental data is in the form of macroscopic  $pK_a$  values, the best way to evaluate computational  
867 predictions is by calculating predicted macroscopic  $pK_a$ . With the conversion of microscopic  $pK_a$  to macroscopic  $pK_a$ s all the  
868 structural information about the titration site is lost and only remaining information is the total charge of macroscopic ion-  
869 ization states. Unfortunately, most macroscopic  $pK_a$  measurements including potentiometric and spectrophotometric meth-  
870 ods do not capture the absolute charge of the macrostates. The spectrophotometric method does not measure charge at all.  
871 The potentiometric method can only capture the relative charge change between macrostates. Only pH-dependent solubility  
872 based  $pK_a$  estimations can differentiate the neutral and charged states from one another. So it is very common to have ex-  
873 perimental datasets of macroscopic  $pK_a$  without any charge or protonation position information regarding the macrostates.  
874 This causes an issue of assigning predicted and experimental  $pK_a$  values before any error statistics can be calculated. As deline-  
875 ated by Fraczkiewicz et. al. the fairest and reasonable solution for  $pK_a$  matching problem involves an assignment algorithm  
876 that preserves the order of predicted and experimental microstates and uses the principle of smallest differences to pair val-  
877 ues [22]. We recommend Hungarian matching with the squared error cost function. The algorithm is available in SciPy package  
878 (`scipy.optimize.linear_sum_assignment`) [30]. In addition to the analysis of numerical error statistics after Hungarian matching,  
879 at the very least number of missing and extra  $pK_a$  predictions must be reported based on unmatched  $pK_a$  values. Missing or  
880 extra  $pK_a$  predictions point to a problem with capturing the right number of ionization states within the pH interval of the exper-  
881 imental measurements. We have demonstrated that for microscopic  $pK_a$  predictions performance analysis based in Hungarian  
882 matching results in overly optimistic and misleading results, instead the employed microstate-based matching provided a more  
883 realistic assessment.

884 We allowed three different submission types in SAMPL6 to capture all the necessary information related to  $pK_a$  predictions.  
885 These were (1) macroscopic  $pK_a$  values, (2) microscopic  $pK_a$  values and microstate pair identities, (3) fractional population of  
886 microstates with respect to pH. We realized later that collecting fractional populations of microstates was redundant since mi-  
887 croscopic  $pK_a$  values and microstate pairs capture all the necessary information to construct fractional population vs. pH curves.  
888 Only microscopic and macroscopic  $pK_a$  values were used for the challenge analysis presented in this paper. While exploring ways  
889 to evaluate SAMPL6  $pK_a$  Challenge results, we developed a better way to capture microscopic  $pK_a$  predictions as presented in  
890 an earlier paper [27]. This alternative reporting format consists of charge and relative free energy of microstates with respect to  
891 a reference microstate and pH. This approach presents the most concise method of capturing all necessary information regard-  
892 ing microscopic  $pK_a$  predictions and allows calculation of predicted microscopic  $pK_a$ s, microstate population with respect to pH,  
893 macroscopic  $pK_a$ s, macroscopic population with respect to pH, and tautomer ratios. Still, there may be methods developed to  
894 predict macroscopic  $pK_a$ s directly instead of computing it from microstate predictions that justifies allowing a macroscopic  $pK_a$   
895 reporting format. In future challenges, we recommend collecting  $pK_a$  predictions with two submission types: (1) macroscopic  
896  $pK_a$  values together with the charges of the macrostates and (2) microstates, their total charge, and relative free energies with  
897 respect to a specified reference microstate and pH. This approach is being used in SAMPL7.

898 In SAMPL6 we provided an enumerated list of microstates and their assigned microstate IDs because we were worried about  
899 parsing submitted microstates in SMILES from different sources correctly. There were two disadvantages to this approach. First,  
900 this list of enumerated microstates was used as input by some participants which was not our intention. Second, the first it-  
901 eration of enumerated microstates was not complete. We had to add new microstates and assign them microstate IDs for a  
902 couple of rounds until reaching a complete list. In future challenges, a better way of handling the problem of capturing predicted  
903 microstates would be asking participants to specify the predicted protonation states themselves and assigning identifiers after  
904 the challenge deadline to aid comparative analysis. This would prevent the partial unblinding of protonation states and allow  
905 the assessment of whether methods can predict all the relevant states independently, without relying on a provided list of mi-

906 crostates. Predicted states can be submitted as mol2 files that represents the microstate with explicit hydrogens. The organizers  
907 must only provide the microstate that was selected as the reference state for the relative microstate free energy calculations.  
908

909 In the SAMPL6  $pK_a$  Challenge there was not a requirement that prediction sets should report predictions for all compounds.  
910 Some participants reported predictions for only a subset of compounds which may have led these methods to look more ac-  
911 curate than others, due to missing predictions. In the future, it will be better to allow submissions of only complete sets for a  
912 better comparison of method performance.

913 A wide range of methods participated in the SAMPL6  $pK_a$  Challenge from very fast QSPR methods to QM methods with a  
914 high-level of theory and extensive exploration of conformational ensembles. In the future, it would be interesting to capture  
915 computing costs in terms of average compute hours per molecule. This can provide guidance to future users of  $pK_a$  prediction  
916 methods for selection of which method to use.

917 Future blind challenges can maximize learning opportunities by evaluating predictions of different physicochemical proper-  
918 ties for the same molecules in consecutive challenges. In SAMPL6 we organized both  $pK_a$  and  $\log P$  challenges. Unfortunately  
919 only a subset of compounds in  $pK_a$  datasets were suitable for the potentiometric  $\log P$  measurements. Still, comparing pre-  
920 diction performance of common compounds in both challenges can lead to beneficial insights especially for physical modeling  
921 techniques if there are common aspects that are beneficial or detrimental to prediction performance. For example, in SAMPL6  
922  $pK_a$  and  $\log P$  Challenges COSMO-RS and EC-RISM solvation models achieved good performance. Having access to a variety  
923 of physicochemical property measurements can also help identification of error sources. For example, dominant microstates  
924 determined for  $pK_a$  challenge can provide information to check if correct tautomers are modeling in a  $\log P$  or  $\log D$  challenge.  
925  $pK_a$  prediction is a requirement for  $\log D$  prediction and experimental  $pK_a$  values can help diagnosing the source of errors in  
926  $\log D$  predictions better. The physical challenges in SAMPL7, which is currently running with a deadline of September 30th, 2020,  
927 follow this principle and include both  $pK_a$ ,  $\log P$ , and membrane permeability properties for a set of monoprotic compounds.  
928 We hope that future  $pK_a$  challenges can focus on multiprotic drug-like compounds with microscopic  $pK_a$  measurements for an  
929 in-depth analysis.

## 929 4 Conclusion

930 The first SAMPL6  $pK_a$  Challenge focused on kinase inhibitor like molecules to assess the performance of  $pK_a$  predictions for  
931 drug-like molecules. With wide participation we had an opportunity to prospectively evaluate  $pK_a$  predictions spanning vari-  
932 ous empirical and QM based approaches. A small number of popular  $pK_a$  prediction methods that were missing from blind  
933 submissions were added as reference calculations after the challenge deadline.

934 The experimental dataset consisted of spectrophotometric measurements of 24 molecules and some of which were multi-  
935 protic. There was also experimental data on the dominant microstate sequence of a subset of the challenge molecules, but  
936 not direct microscopic  $pK_a$  measurements. We have performed a comparative analysis of methods represented in the blind  
937 challenge in terms of both macroscopic and microscopic  $pK_a$  prediction performance avoiding any assumptions about the ex-  
938 perimental  $pK_a$ s.

939 Here, we used Hungarian matching to assign predicted and experimental values for the calculation of accuracy and corre-  
940 lation statistics, because the majority of the experimental data was macroscopic  $pK_a$  values. In addition to evaluating error in  
941 predicted  $pK_a$  values, we also reported the macroscopic  $pK_a$  errors that were not captured by the match between experimental  
942 and predicted  $pK_a$  values. These were extra or missing  $pK_a$  predictions which are important indicators that predictions are failing  
943 to capture the correct ionization states.

944 We evaluated microscopic  $pK_a$  predictions utilizing the experimental dominant microstate sequence data of eight molecules.  
945 This experimental data allowed us to use microstate-based matching for evaluating the accuracy of microscopic  $pK_a$  values  
946 in a more realistic way. We have determined that QM and LFER predictions had lower accuracy in determining the dominant  
947 tautomer of the charged microstates than the neutral states. For both macroscopic and microscopic  $pK_a$  predictions we have  
948 determined methods that were consistently well-performing according to multiple statistical metrics. Focusing on the com-  
949 parison of molecules instead of methods for macroscopic  $pK_a$  prediction accuracy indicated molecules with sulfur-containing  
950 heterocycles, iodo, and bromo groups suffered from lower  $pK_a$  prediction accuracy.

951 The overall performance of  $pK_a$  predictions as captured in this challenge is concerning for the application of  $pK_a$  prediction  
952 methods in computer-aided drug design. Many computational methods for predicting target affinities and physicochemical  
953 properties rely on  $pK_a$  predictions for determining relevant protonation states and the free energy penalty of such states. 1 unit  
954 of  $pK_a$  error is an optimistic estimate of current macroscopic  $pK_a$  predictions for drug-like molecules based on SAMPL6 Challenge

where errors in predicting the correct number of ionization states or determining the correct dominant microstate were also common to many methods. In the absence of other sources of errors, we showed that 1 unit over- or underestimation of the  $pK_a$  of a ligand can cause significant errors in the overall binding affinity calculation due to errors in multiple protonation state correction factor.

The SAMPL6 GitHub Repository contains all information regarding the challenge structure, experimental data, blind prediction submission sets, and evaluation of methods. The repository will be useful for future follow up analysis and the experimental measurements can continue to serve as a benchmark dataset for testing methods.

In this article, we aimed to demonstrate not only the comparative analysis of the  $pK_a$  prediction performance of contemporary methods for drug-like molecules, but also to propose a stringent  $pK_a$  prediction evaluation strategy that takes into account differences in microscopic and macroscopic  $pK_a$  definitions. We hope that this study will guide and motivate further improvement of  $pK_a$  prediction methods.

## 5 Code and data availability

- SAMPL6  $pK_a$  challenge instructions, submissions, experimental data and analysis is available at <https://github.com/samplchallenges/SAMPL6>

## 6 Overview of supplementary information

Contents of the Supplementary Information:

- TABLE S1: SMILES and InChI identifiers of SAMPL6  $pK_a$  Challenge molecules.
- TABLE S2: Evaluation statistics calculated for all macroscopic  $pK_a$  prediction submissions based on Hungarian match for 24 molecules.
- TABLE S3: Evaluation statistics calculated for all microscopic  $pK_a$  prediction submissions based on Hungarian match for 8 molecules with NMR data.
- TABLE S4: Evaluation statistics calculated for all microscopic  $pK_a$  prediction submissions based on microstate match for 8 molecules with NMR data.
- FIGURE S1: Dominant microstates of 8 molecules were determined based on NMR measurements.
- FIGURE S2: MAE of macroscopic  $pK_a$  predictions of each molecule did not show any significant correlation with any molecular descriptor.
- FIGURE S3: The value of macroscopic  $pK_a$  was not a factor affecting prediction error seen in SAMPL6 Challenge according to the analysis with Hungarian matching.
- FIGURE S4: There was low agreement between experimental dominant microstate pairs and the predicted microstate pairs selected by Hungarian algorithm for microscopic  $pK_a$  predictions.

Extra files included in *SAMPL6-supplementary-documents.tar.gz*:

- SAMPL6-pKa-chemical-identifiers-table.csv
- macroscopic-pKa-statistics-24mol-hungarian-match.csv
- microscopic-pKa-statistics-8mol-hungarian-match-table.csv
- microscopic-pKa-statistics-8mol-microstate-match-table.csv
- experimental-microstates-of-8mol-based-on-NMR.csv
- enumerate-microstates-with-Epik-and-OpenEye-QUACPAC.ipynb
- molecule\_ID\_and\_SMILES.csv

## 7 Author Contributions

Conceptualization, MI, JDC ; Methodology, MI, JDC, ASR ; Software, MI, AR, ASR ; Formal Analysis, MI, ASR ; Investigation, MI ; Resources, JDC, DLM; Data Curation, MI ; Writing-Original Draft, MI; Writing - Review and Editing, MI, JDC, ASR, AR, DLM; Visualization, MI, AR ; Supervision, JDC, DLM ; Project Administration, MI ; Funding Acquisition, JDC, DLM.

## 8 Acknowledgments

We would like to acknowledge the infrastructure and website support of Mike Chiu that allowed a seamless collection of challenge submissions. Mike Chiu also provided assistance with constructing a submission validation script to ensure all submissions

999 adhered to the machine-readable format. We are grateful to Kiril Lanevskij for suggesting the Hungarian algorithm for matching  
1000 experimental and predicted  $pK_a$  values. We would like to thank Thomas Fox for providing MoKa reference calculations. We  
1001 acknowledge Caitlin Bannan for guidance on defining a working microstate definition for the challenge and guidance for design-  
1002 ing the challenge. We thank Brad Sherborne for his valuable insights at the conception of the  $pK_a$  challenge and connecting us  
1003 with Timothy Rhodes and Dorothy Levorse who were able to provide resources and expertise for experimental measurements  
1004 performed at MRL. We acknowledge Paul Czodrowski who provided feedback on multiple stages of this work: challenge construc-  
1005 tion, purchasable compound selection, and manuscript draft. MI, JDC, and DLM gratefully acknowledge support from NIH grant  
1006 R01GM124270 supporting the SAMPL Blind Challenges. MI, ASR, AR, and JDC acknowledge support from the Sloan Kettering  
1007 Institute. JDC acknowledges support from NIH grant P30 CA008748. DLM appreciates financial support from the National Insti-  
1008 tutes of Health (1R01GM108889-01) and the National Science Foundation (CHE 1352608). MI acknowledges Doris J. Hutchinson  
1009 Fellowship. MI, ASR, AR, and JDC are grateful to OpenEye Scientific for providing a free academic software license for use in this  
1010 work. MI, ASR, AR, and JDC thank Janos Fejervari and ChemAxon team that gave us permission to include ChemAxon/Chemicalize  
1011  $pK_a$  predictions as a reference prediction in challenge analysis.

## 1012 9 Disclaimers

1013 The content is solely the responsibility of the authors and does not necessarily represent the official views of the National  
1014 Institutes of Health.

## 1015 10 Disclosures

1016 JDC was a member of the Scientific Advisory Board for Schrödinger, LLC during part of this study. JDC and DLM are current  
1017 members of the Scientific Advisory Board of OpenEye Scientific Software, and DLM is an Open Science Fellow with Silicon Ther-  
1018 apeutics. The Chodera laboratory receives or has received funding from multiple sources, including the National Institutes of  
1019 Health, the National Science Foundation, the Parker Institute for Cancer Immunotherapy, Relay Therapeutics, Entasis Therapeu-  
1020 tics, Vir Biotechnology, Silicon Therapeutics, EMD Serono (Merck KGaA), AstraZeneca, Vir Biotechnology, XtalPi, the Molecular  
1021 Sciences Software Institute, the Starr Cancer Consortium, the Open Force Field Consortium, Cycle for Survival, a Louis V. Ger-  
1022 stner Young Investigator Award, The Einstein Foundation, and the Sloan Kettering Institute. A complete list of funding can be  
1023 found at <http://choderlab.org/funding>.

## 1024 References

- 1025 [1] Manallack DT, Prankerd RJ, Yuriev E, Oprea TI, Chalmers DK. The Significance of Acid/Base Properties in Drug Discovery. *Chem Soc Rev.* 2013; 42(2):485–496. doi: [10.1039/C2CS35348B](https://doi.org/10.1039/C2CS35348B).
- 1026 [2] Manallack DT, Prankerd RJ, Nassta GC, Ursu O, Oprea TI, Chalmers DK. A Chemogenomic Analysis of Ionization Constants-Implications for Drug Discovery. *ChemMedChem.* 2013 Feb; 8(2):242–255. doi: [10.1002/cmdc.201200507](https://doi.org/10.1002/cmdc.201200507).
- 1027 [3] de Oliveira C, Yu HS, Chen W, Abel R, Wang L. Rigorous Free Energy Perturbation Approach to Estimating Relative Binding Affinities  
1028 between Ligands with Multiple Protonation and Tautomeric States. *Journal of Chemical Theory and Computation.* 2019 Jan; 15(1):424–435.  
doi: [10.1021/acs.jctc.8b00826](https://doi.org/10.1021/acs.jctc.8b00826).
- 1029 [4] Darvey IG. The Assignment of  $pK_a$  Values to Functional Groups in Amino Acids. *Biochemical Education.* 1995 Apr; 23(2):80–82. doi:  
1030 [10.1016/0307-4412\(94\)00150-N](https://doi.org/10.1016/0307-4412(94)00150-N).
- 1031 [5] Bodner GM. Assigning the  $pK_a$ 's of Polyprotic Acids. *Journal of Chemical Education.* 1986 Mar; 63(3):246. doi: [10.1021/ed063p246](https://doi.org/10.1021/ed063p246).
- 1032 [6] Murray R. Microscopic Equilibria. *Analytical Chemistry.* 1995 Aug; p. 1.
- 1033 [7] Işık M, Levorse D, Rustenburg AS, Ndukwe IE, Wang H, Wang X, Reibarkh M, Martin GE, Makarov AA, Mobley DL, Rhodes T, Chodera JD.  $pK_a$  Measurements for the SAMPL6 Prediction Challenge for a Set of Kinase Inhibitor-like Fragments. *Journal of Computer-Aided Molecular  
1034 Design.* 2018 Oct; 32(10):1117–1138. doi: [10.1007/s10822-018-0168-0](https://doi.org/10.1007/s10822-018-0168-0).
- 1035 [8] Bochevarov AD, Watson MA, Greenwood JR, Philipp DM. Multiconformation, Density Functional Theory-Based  $pK_a$  Prediction in Application to Large, Flexible Organic Molecules with Diverse Functional Groups. *Journal of Chemical Theory and Computation.* 2016 Dec;  
1036 12(12):6001–6019. doi: [10.1021/acs.jctc.6b00805](https://doi.org/10.1021/acs.jctc.6b00805).
- 1037 [9] Selwa E, Kenney IM, Beckstein O, Iorga BI. SAMPL6: Calculation of Macroscopic  $pK_a$  Values from Ab Initio Quantum Mechanical Free  
1038 Energies. *Journal of Computer-Aided Molecular Design.* 2018 Oct; 32(10):1203–1216. doi: [10.1007/s10822-018-0138-6](https://doi.org/10.1007/s10822-018-0138-6).

- 1044 [10] **Pickard FC**, König G, Tofoleanu F, Lee J, Simonett AC, Shao Y, Ponder JW, Brooks BR. Blind Prediction of Distribution in the SAMPL5  
1045 Challenge with QM Based Protomer and pK<sub>a</sub> Corrections. *Journal of Computer-Aided Molecular Design*. 2016 Nov; 30(11):1087–1100. doi:  
1046 [10.1007/s10822-016-9955-7](https://doi.org/10.1007/s10822-016-9955-7).
- 1047 [11] **Bannan CC**, Mobley DL, Skillman AG. SAMPL6 Challenge Results from \$\$pK\_a\$\$ Predictions Based on a General Gaussian Process Model.  
1048 *Journal of Computer-Aided Molecular Design*. 2018 Oct; 32(10):1165–1177. doi: [10.1007/s10822-018-0169-z](https://doi.org/10.1007/s10822-018-0169-z).
- 1049 [12] **Işık M**, Levorse D, Mobley DL, Rhodes T, Chodera JD. Octanol–Water Partition Coefficient Measurements for the SAMPL6 Blind Prediction  
1050 Challenge. *Journal of Computer-Aided Molecular Design*. 2020 Apr; 34(4):405–420. doi: [10.1007/s10822-019-00271-3](https://doi.org/10.1007/s10822-019-00271-3).
- 1051 [13] **Işık M**, Bergazin TD, Fox T, Rizzi A, Chodera JD, Mobley DL. Assessing the Accuracy of Octanol–Water Partition Coefficient Predictions in the  
1052 SAMPL6 Part II Log P Challenge. *Journal of Computer-Aided Molecular Design*. 2020 Apr; 34(4):335–370. doi: [10.1007/s10822-020-00295-0](https://doi.org/10.1007/s10822-020-00295-0).
- 1053 [14] **Kogej T**, Muresan S. Database Mining for pKa Prediction. *Current Drug Discovery Technologies*. 2005; 2(4):221–229. doi:  
1054 [10.2174/157016305775202964](https://doi.org/10.2174/157016305775202964).
- 1055 [15] **Perrin DD**, Dempsey B, Serjeant EP. pKa Prediction for Organic Acids and Bases. 1 ed. London and New York: Chapman and Hall; 1981.
- 1056 [16] **Hammett LP**. Physical Organic Chemistry. New York: McGraw-Hill; 1940.
- 1057 [17] **Taft RW**, Lewis IC. Evaluation of Resonance Effects on Reactivity by Application of the Linear Inductive Energy Relationship. V. Concerning  
1058 a  $\sigma R$  Scale of Resonance Effects<sup>1,2</sup>. *Journal of the American Chemical Society*. 1959; 81(20):5343–5352. doi: [10.1021/ja01529a025](https://doi.org/10.1021/ja01529a025).
- 1059 [18] **Xing L**, Glen RC, Clark RD. Predicting  $p K_a$  by Molecular Tree Structured Fingerprints and PLS. *Journal of Chemical Information and  
1060 Computer Sciences*. 2003 May; 43(3):870–879. doi: [10.1021/ci020386s](https://doi.org/10.1021/ci020386s).
- 1061 [19] **Zhang J**, Kleinöder T, Gasteiger J. Prediction of  $p K_a$  Values for Aliphatic Carboxylic Acids and Alcohols with Empirical Atomic Charge  
1062 Descriptors. *Journal of Chemical Information and Modeling*. 2006 Nov; 46(6):2256–2266. doi: [10.1021/ci060129d](https://doi.org/10.1021/ci060129d).
- 1063 [20] **Cruciani G**, Milletti F, Storchi L, Sforza G, Goracci L. *In Silico*  $p K_a$  Prediction and ADME Profiling. *Chemistry & Biodiversity*. 2009 Nov;  
1064 6(11):1812–1821. doi: [10.1002/cbdv.200900153](https://doi.org/10.1002/cbdv.200900153).
- 1065 [21] **Milletti F**, Storchi L, Sforza G, Cruciani G. New and Original  $p K_a$  Prediction Method Using Grid Molecular Interaction Fields. *Journal of  
1066 Chemical Information and Modeling*. 2007 Nov; 47(6):2172–2181. doi: [10.1021/ci700018y](https://doi.org/10.1021/ci700018y).
- 1067 [22] **Fraczkiewicz R**. In Silico Prediction of Ionization. In: *Reference Module in Chemistry, Molecular Sciences and Chemical Engineering*; Elsevier;  
1068 2013. doi: [10.1016/B978-0-12-409547-2.02610-X](https://doi.org/10.1016/B978-0-12-409547-2.02610-X).
- 1069 [23] Simulations Plus ADMET Predictor v8.5;. Simulations Plus, Lancaster, CA, 2018. <https://www.simulations-plus.com/software/admetpredictor/physicochemical-biopharmaceutical/>.
- 1070 [24] Special Issue: SAMPL6 (Statistical Assessment of the Modeling of Proteins and Ligands); October 2018. Volume 32, Issue 10. *Journal of  
1071 Computer-Aided Molecular Design*.
- 1072 [25] **Shelley JC**, Cholleti A, Frye LL, Greenwood JR, Timlin MR, Uchimaya M. Epik: A Software Program for pK<sub>a</sub> Prediction and Protonation State  
1073 Generation for Drug-like Molecules. *Journal of Computer-Aided Molecular Design*. 2007 Dec; 21(12):681–691. doi: [10.1007/s10822-007-9133-z](https://doi.org/10.1007/s10822-007-9133-z).
- 1074 [26] QUACPAC Toolkit 2017.Feb.1;. OpenEye Scientific Software, Santa Fe, NM. <http://www.eyesopen.com>.
- 1075 [27] **Gunner MR**, Murakami T, Rustenburg AS, Işık M, Chodera JD. Standard State Free Energies, Not pK<sub>a</sub>s, Are Ideal for Describing Small  
1076 Molecule Protonation and Tautomeric States. *Journal of Computer-Aided Molecular Design*. 2020 May; 34(5):561–573. doi: [10.1007/s10822-020-00280-7](https://doi.org/10.1007/s10822-020-00280-7).
- 1077 [28] **Kuhn HW**. The Hungarian Method for the Assignment Problem. *Naval Research Logistics Quarterly*. 1955 Mar; 2(1-2):83–97. doi:  
1081 [10.1002/nav.3800020109](https://doi.org/10.1002/nav.3800020109).
- 1082 [29] **Munkres J**. Algorithms for the Assignment and Transportation Problems. *J SIAM*. 1957 Mar; 5(1):32–28.
- 1083 [30] SciPy v1.3.1, Linear Sum Assignment Documentation; Sep 27, 2019. The SciPy community. [https://docs.scipy.org/doc/scipy-1.3.1/reference/generated/scipy.optimize.linear\\_sum\\_assignment.html](https://docs.scipy.org/doc/scipy-1.3.1/reference/generated/scipy.optimize.linear_sum_assignment.html).
- 1084 [31] OpenEye pKa Prospector;. OpenEye Scientific Software, Santa Fe, NM. Accessed on Jan 23, 2018. <https://www.eyesopen.com/pka-prospector>.
- 1085 [32] ACD/pKa GALAS (ACD/Percepta Kernel v1.6);. Advanced Chemistry Development, Inc., Toronto, ON, Canada, 2018. <https://www.acdlabs.com/products/percepta/predictors/pKa/>.

- 1088 [33] ACD/pKa Classic (ACD/Percepta Kernel v1.6);. Advanced Chemistry Development, Inc., Toronto, ON, Canada, 2018. <https://www.acdlabs.com/products/percepta/predictors/pKa/>.
- 1089
- 1090 [34] Chemicalize v18.23 (ChemAxon MarvinSketch v18.23);. ChemAxon, Budapest, Hungary, 2018. <https://docs.chemaxon.com/display/docs/pKa+Plugin>.
- 1091
- 1092 [35] MoKa;. Molecular Discovery, Hertfordshire, UK, 2018. <https://www.moldiscovery.com/software/moka/>.
- 1093 [36] Zeng Q, Jones MR, Brooks BR. Absolute and Relative pKa Predictions via a DFT Approach Applied to the SAMPL6 Blind Challenge. *Journal of Computer-Aided Molecular Design*. 2018 Oct; 32(10):1179–1189. doi: [10.1007/s10822-018-0150-x](https://doi.org/10.1007/s10822-018-0150-x).
- 1094
- 1095 [37] Bochevarov AD, Harder E, Hughes TF, Greenwood JR, Braden DA, Philipp DM, Rinaldo D, Halls MD, Zhang J, Friesner RA. Jaguar: A High-  
1096 Performance Quantum Chemistry Software Program with Strengths in Life and Materials Sciences. *International Journal of Quantum  
1097 Chemistry*. 2013 Sep; 113(18):2110–2142. doi: [10.1002/qua.24481](https://doi.org/10.1002/qua.24481).
- 1098 [38] Tielker N, Eberlein L, Güssregen S, Kast SM. The SAMPL6 Challenge on Predicting Aqueous pKa Values from EC-RISM Theory. *Journal of  
1099 Computer-Aided Molecular Design*. 2018 Oct; 32(10):1151–1163. doi: [10.1007/s10822-018-0140-z](https://doi.org/10.1007/s10822-018-0140-z).
- 1100 [39] Klamt A, Eckert F, Diedenhofen M, Beck ME. First Principles Calculations of Aqueous  $pK_a$  Values for Organic and Inorganic Acids Using  
1101 COSMO-RS Reveal an Inconsistency in the Slope of the  $pK_a$  Scale. *The Journal of Physical Chemistry A*. 2003 Nov; 107(44):9380–9386. doi:  
1102 [10.1021/jp034688o](https://doi.org/10.1021/jp034688o).
- 1103 [40] Eckert F, Klamt A. Accurate Prediction of Basicity in Aqueous Solution with COSMO-RS. *Journal of Computational Chemistry*. 2006 Jan;  
1104 27(1):11–19. doi: [10.1002/jcc.20309](https://doi.org/10.1002/jcc.20309).
- 1105 [41] Pracht P, Wilcken R, Udvarhelyi A, Rodde S, Grimme S. High Accuracy Quantum-Chemistry-Based Calculation and Blind Prediction of  
1106 Macroscopic pKa Values in the Context of the SAMPL6 Challenge. *Journal of Computer-Aided Molecular Design*. 2018 Oct; 32(10):1139–  
1107 1149. doi: [10.1007/s10822-018-0145-7](https://doi.org/10.1007/s10822-018-0145-7).
- 1108 [42] Prasad S, Huang J, Zeng Q, Brooks BR. An Explicit-Solvent Hybrid QM and MM Approach for Predicting pKa of Small Molecules in SAMPL6  
1109 Challenge. *Journal of Computer-Aided Molecular Design*. 2018 Oct; 32(10):1191–1201. doi: [10.1007/s10822-018-0167-1](https://doi.org/10.1007/s10822-018-0167-1).
- 1110 [43] Robert Fraczkiewicz MW, SAMPL6 pKa Challenge: Predictions of ionization constants performed by the S+pKa method implemented in  
1111 ADMET Predictor software; February 22, 2018. The Joint D3R/SAMPL Workshop 2018. <https://drugdesigndata.org/about/d3r-2018-workshop>.
- 1112 [44] Mansouri K, Cariello NF, Korotcov A, Tkachenko V, Grulke CM, Sprankle CS, Allen D, Casey WM, Kleinstreuer NC, Williams AJ. Open-  
1113 Source QSAR Models for pKa Prediction Using Multiple Machine Learning Approaches. *Journal of Cheminformatics*. 2019 Dec; 11(1). doi:  
1114 [10.1186/s13321-019-0384-1](https://doi.org/10.1186/s13321-019-0384-1).
- 1115 [45] Baltruschat M, Czodrowski P. Machine Learning Meets pKa [Version 2; Peer Review: 2 Approved]. *F1000Research*. 2020; 9 (Chem Inf  
1116 Sci)(113). doi: [10.12688/f1000research.22090.2](https://doi.org/10.12688/f1000research.22090.2).
- 1117 [46] Hunt P, Hosseini-Gerami L, Chrien T, Plante J, Ponting DJ, Segall M. Predicting  $pK_a$  Using a Combination of Semi-Empirical Quan-  
1118 tum Mechanics and Radial Basis Function Methods. *Journal of Chemical Information and Modeling*. 2020 Jun; 60(6):2989–2997. doi:  
1119 [10.1021/acs.jcim.0c00105](https://doi.org/10.1021/acs.jcim.0c00105).
- 1120 [47] Zdrazil B, Guha R. The Rise and Fall of a Scaffold: A Trend Analysis of Scaffolds in the Medicinal Chemistry Literature. *Journal of Medicinal  
1121 Chemistry*. 2018 Jun; 61(11):4688–4703. doi: [10.1021/acs.jmedchem.7b00954](https://doi.org/10.1021/acs.jmedchem.7b00954).
- 1122 [48] Ertl P, Altmann E, McKenna JM. The Most Common Functional Groups in Bioactive Molecules and How Their Popularity Has Evolved over  
1123 Time. *Journal of Medicinal Chemistry*. 2020 Aug; 63(15):8408–8418. doi: [10.1021/acs.jmedchem.0c00754](https://doi.org/10.1021/acs.jmedchem.0c00754).
- 1124 [49] OEMolProp Toolkit 2017.Feb.1;. OpenEye Scientific Software, Santa Fe, NM. <http://www.eyesopen.com>.

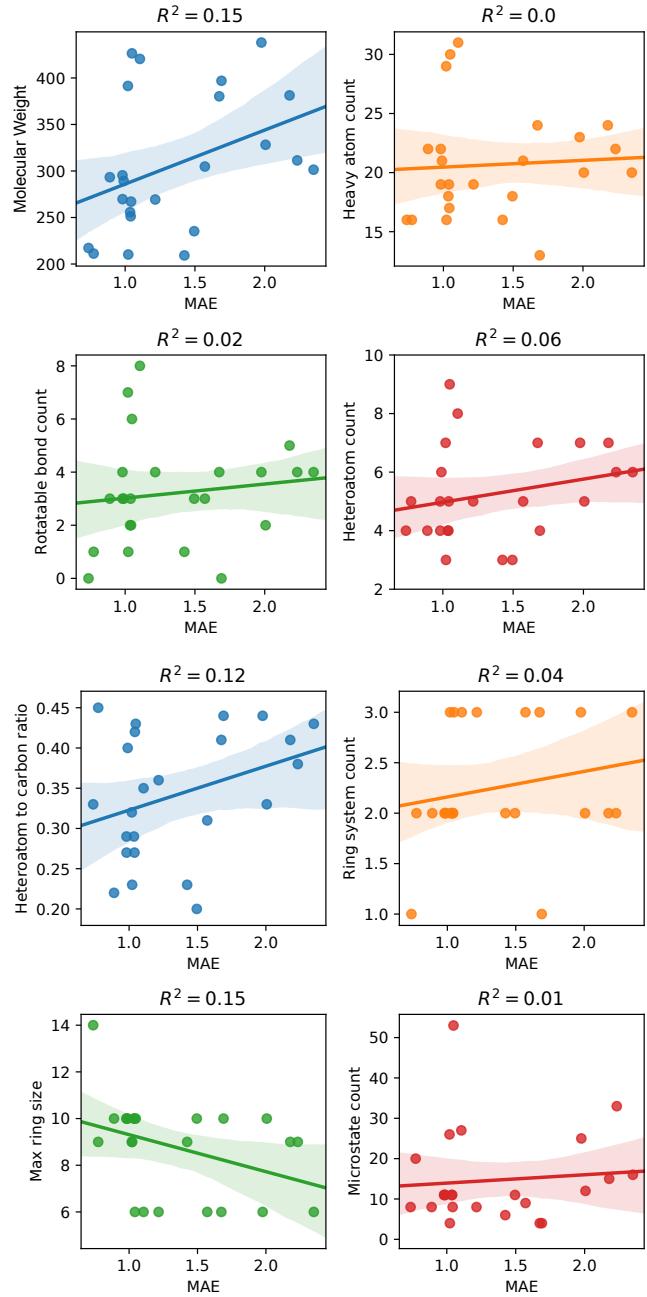
**Table S1. SMILES and InChI identifiers of SAMPL6 pK<sub>a</sub> Challenge molecules.** A CSV version of this table can be found in *SAMPL6-supplementary-documents.tar.gz*.

SAMPL6 Molecule ID	Isomeric SMILES	InChI
SM01	c1cc2c(cc1O)c3c(o2)C(=O)NCCC3	InChI=1S/C12H11NO3/c14-7-3-4-10-9(6-7)8-2-1-5-13-12(15)11(8)16-10/h3-4,6,14H,1-2,5H2,(H,13,15)
SM02	c1ccc2c(c1)c(ncn2)Nc3cccc(c3)C(F)(F)	InChI=1S/C15H10F3N3/c16-15(17,18)10-4-3-5-11(8-10)21-14-12-6-1-2-7-13(12)19-9-20-14/h1-9H,(H,19,20,21)
SM03	c1ccc(cc1)Cc2nnnc(s2)NC(=O)c3cccs3	InChI=1S/C14H11N3OS2/c18-13(11-7-4-8-19-11)15-14-17-16-12(20-14)9-10-5-2-1-3-6-10/h1-8H,9H2,(H,15,17,18)
SM04	c1ccc2c(c1)c(ncn2)NCc3ccc(cc3)Cl	InChI=1S/C15H12ClN3/c16-12-7-5-11(6-8-12)9-17-15-13-3-1-2-4-14(13)18-10-19-15/h1-8,10H,9H2,(H,17,18,19)
SM05	c1ccc(c(c1)NC(=O)c2ccc(o2)Cl)N3CCCCC3	InChI=1S/C16H17ClN2O2/c17-15-9-8-14(21-15)16(20)18-12-6-2-3-7-13(12)19-10-4-1-5-11-19/h2-3,6-9H,1,4-5,10-11H2,(H,18,20)
SM06	c1cc2ccnc2c(c1)NC(=O)c3cc(cnc3)Br	InChI=1S/C15H10BrN3O/c16-12-7-11(8-17-9-12)15(20)19-13-5-1-3-10-4-2-6-18-14(10)13/h1-9H,(H,19,20)
SM07	c1ccc(cc1)CNc2c3cccc3ncn2	InChI=1S/C15H13N3/c1-2-6-12(7-3-1)10-16-15-13-8-4-5-9-14(13)17-11-18-15/h1-9,11H,10H2,(H,16,17,18)
SM08	Cc1ccc2c(c1)c(c(c(=O)[nH]2)CC(=O)O)c3cccc3	InChI=1S/C18H15NO3/c1-11-7-8-15-13(9-11)17(12-5-3-2-4-6-12)14(10-16(20)21)18(22)19-15/h2-9H,10H2,1H3,(H,19,22)(H,20,21)
SM09	COc1cccc(c1)Nc2c3cccc3ncn2.Cl	InChI=1S/C15H13N3O.CIH/c1-19-12-6-4-5-11(9-12)18-15-13-7-2-3-8-14(13)16-10-17-15;/h2-10H,1H3,(H,16,17,18);1H
SM10	c1ccc(cc1)C(=O)NCC(=O)Nc2nc3cccc3s2	InChI=1S/C16H13N3O2S/c20-14(10-17-15(21)11-6-2-1-3-7-11)19-16-18-1-2-8-4-5-9-13(12)22-16/h1-9H,10H2,(H,17,21)(H,18,19,20)
SM11	c1ccc(cc1)n2c3c(cn2)c(ncn3)N	InChI=1S/C11H9N5/c12-10-9-6-15-16(11(9)14-7-13-10)8-4-2-1-3-5-8/h1-7H,(H,2,12,13,14)
SM12	c1ccc2c(c1)c(ncn2)Nc3cccc(c3)Cl.Cl	InChI=1S/C14H10ClN3.CIH/c15-10-4-3-5-11(8-10)18-14-12-6-1-2-7-13(12)16-9-17-14;/h1-9H,(H,16,17,18);1H
SM13	Cc1cccc(c1)Nc2c3cc(c(c3ncn2)OC)OC	InChI=1S/C17H17N3O2/c1-11-5-4-6-12(7-11)20-17-13-8-15(21-2)16(22-3)9-14(13)18-10-19-17/h4-10H,1-3H3,(H,18,19,20)
SM14	c1ccc(cc1)n2ncn3c2ccc(c3)N	InChI=1S/C13H11N3/c14-10-6-7-13-12(8-10)15-9-16(13)11-4-2-1-3-5-11/h1-9H,14H2
SM15	c1ccc2c(c1)ncn2c3ccc(cc3)O	InChI=1S/C13H10N2O/c16-11-7-5-10(6-8-11)15-9-14-12-3-1-2-4-13(12)15/h1-9,16H
SM16	c1cc(c(c(c1)Cl)C(=O)Nc2ccncc2)Cl	InChI=1S/C12H8Cl2N2O/c13-9-2-1-3-10(14)11(9)12(17)16-8-4-6-15-7-5-8/h1-7H,(H,15,16,17)
SM17	c1ccc(cc1)CSc2nnc(o2)c3ccncc3	InChI=1S/C14H11N3OS/c1-2-4-11(5-3-1)10-19-14-17-16-13(18-14)12-6-8-15-9-7-12/h1-9H,10H2
SM18	c1ccc2c(c1)c(=O)[nH]c(n2)CCC(=O)Nc3ncc(s3)Cc4ccc(c(c4)F)F	InChI=1S/C21H16F2N4O2S/c22-15-6-5-12(10-16(15)23)9-13-11-24-21(30-13)27-19(28)8-7-18-25-17-4-2-1-3-14(17)20(29)26-18/h1-6,10-11H,7-9H2,(H,24,27,28)(H,25,26,29)
SM19	CCOc1ccc2c(c1)sc(n2)NC(=O)Cc3ccc(c(c3)Cl)Cl	InChI=1S/C17H14Cl2N2O2S/c1-2-23-11-4-6-14-15(9-11)24-17(20-14)21-6(22)8-10-3-5-12(18)13(9)7-10/h3-7,9H,2,8H2,1H3,(H,20,21,22)
SM20	c1cc(cc(c1)OCc2ccc(cc2Cl)Cl)/C=C/3\C(=O)NC(=O)S3	InChI=1S/C17H11Cl2NO3S/c18-12-5-4-11(14(19)8-12)9-23-13-3-1-2-10(6-13)7-15-16(21)20-17(22)24-15/h1-8H,9H2,(H,20,21,22)/b15-7+
SM21	c1cc(cc(c1)Br)Nc2c(cnc(n2)Nc3cccc(c3)Br)F	InChI=1S/C16H11Br2FN4/c17-10-3-1-5-12(7-10)21-15-14(19)9-20-16(23-15)22-13-6-2-4-11(18)8-13/h1-9H,(H,20,21,22,23)
SM22	c1cc2c(cc(c(c2nc1)O))l	InChI=1S/C9H5l2NO/c10-6-4-7(11)9(13)8-5(6)2-1-3-12-8/h1-4,13H
SM23	CCOC(=O)c1ccc(cc1)Nc2cc(cnc(n2)Nc3ccc(cc3)C(=O)OCC)C	InChI=1S/C23H24N4O4/c1-4-30-21(28)16-6-10-18(11-7-16)25-20-14-15(3)24-23(27-20)26-19-12-8-17(9-13-19)22(29)31-5-2/h6-14H,4-5H2,1-3H3,(H2,24,25,26,27)
SM24	COc1ccc(cc1)c2c3c(ncn3oc2c4ccc(cc4)OC)NCCO	InChI=1S/C22H21N3O4/c1-27-16-7-3-14(4-8-16)18-19-21(23-11-12-26)24-13-25-22(19)29-20(18)15-5-9-17(28-2)10-6-15/h3-10,13,26H,11-12H2,1-2H3,(H,23,24,25)

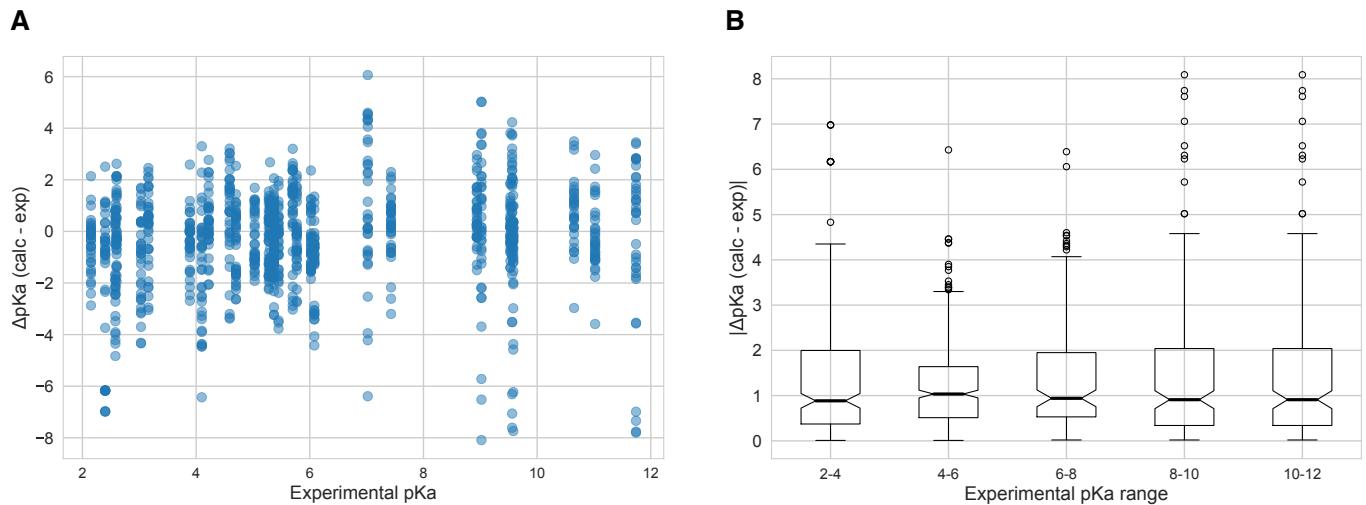
## 11 Supplementary Information

Microstate ID of Deprotonated State (A)	Microstate ID of Protonated State (HA)	Molecule ID	pKa (exp)	pKa SEM (exp)	pKa ID	Microstate identification source
		SM07	6.08	0.01	SM07_pKa1	NMR measurement
		SM14	5.3	0.01	SM14_pKa2	NMR measurement
		SM14	2.58	0.01	SM14_pKa1	NMR measurement
		SM02	5.03	0.01	SM02_pKa1	Estimated based on SM07 NMR measurement
		SM04	6.02	0.01	SM04_pKa1	Estimated based on SM07 NMR measurement
		SM09	5.37	0.01	SM09_pKa1	Estimated based on SM07 NMR measurement
		SM12	5.28	0.01	SM12_pKa1	Estimated based on SM07 NMR measurement
		SM13	5.77	0.01	SM13_pKa1	Estimated based on SM07 NMR measurement
		SM15	8.94	0.01	SM15_pKa2	Estimated based on SM14 NMR measurement
		SM15	4.7	0.01	SM15_pKa1	Estimated based on SM14 NMR measurement

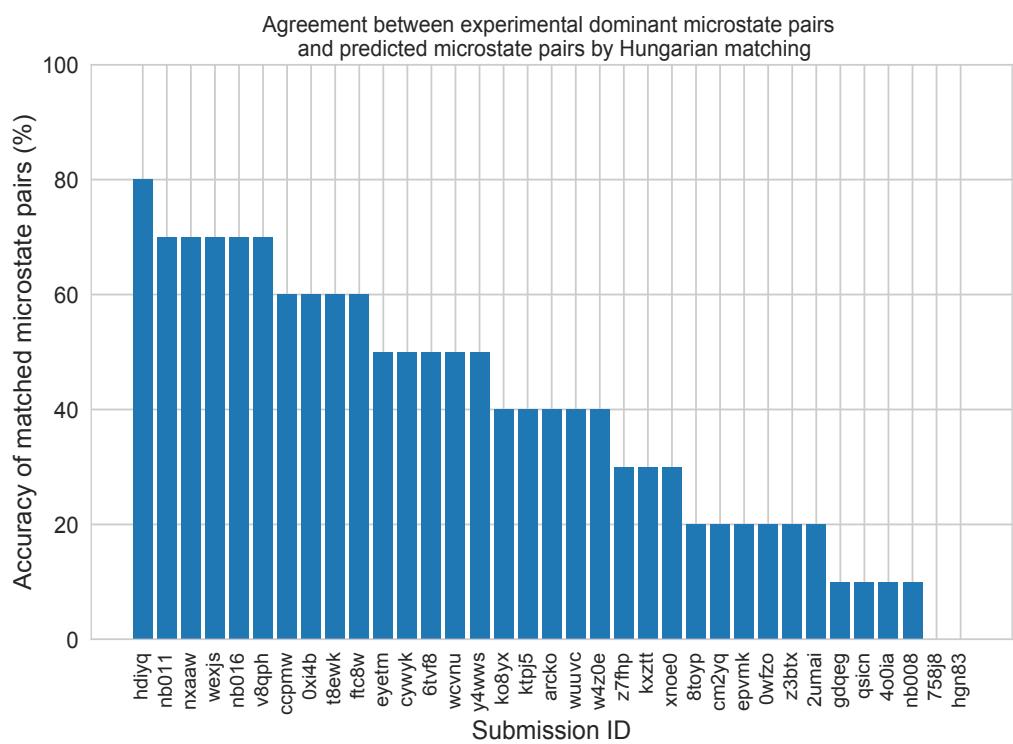
**Figure S1. Dominant microstates of 8 molecules were determined based on NMR measurements.** Dominant microstate sequence of 6 derivatives were determined taking SM07 and SM14 as reference. Matched experimental pK<sub>a</sub> values were determined by spectrophotometric pK<sub>a</sub> measurements [7]. A CSV version of this table can be found in SAMPL6-supplementary-documents.tar.gz.



**Figure S2. MAE of macroscopic  $pK_a$  predictions of each molecule did not show any significant correlation with any molecular descriptor.**  
 Plots show regression lines, 96% confidence intervals of the regression lines, and  $R_2$ . The following molecular descriptors were calculated using OpenEye OEMolProp Toolkit [49].



**Figure S3. The value of macroscopic  $pK_a$ s was not a factor affecting prediction error seen in SAMPL6 Challenge according to the analysis with Hungarian matching.** There was not clear trend between  $pK_a$  prediction error and the true  $pK_a$  error. Very high and very low  $pK_a$  values have similar inaccuracy compared to  $pK_a$  values close to 7. **A** Scatter plot of macroscopic  $pK_a$  prediction error calculated with Hungarian matching vs. experimental  $pK_a$  value **B** Box plot of absolute error of macroscopic  $pK_a$  predictions binned into 2  $pK_a$  unit intervals of experimental  $pK_a$ .



**Figure S4. There was low agreement between experimental dominant microstate pairs and the predicted microstate pairs selected by Hungarian algorithm for microscopic  $pK_a$  predictions.** This analysis could only be performed for 8 molecules with NMR data. Hungarian matching algorithm which matches predicted and experimental values considering only the closeness of the numerical value of  $pK_a$  and it often leads to predicted  $pK_a$  matches that described a different microstates pair than the experimentally observed dominant microstates..

**Table S2. Evaluation statistics calculated for all macroscopic pK<sub>a</sub> prediction submissions based on Hungarian match for 24 molecules.** Methods are represented via their SAMPL6 submission IDs which can be cross-referenced with Table 1 for method details. There are eight error metrics reported: the root-mean-squared error (RMSE), mean absolute error (MAE), mean (signed) error (ME), coefficient of determination ( $R^2$ ), linear regression slope (m), Kendall's Rank Correlation Coefficient ( $\tau$ ), unmatched experimental pK<sub>a</sub>s (number of missing pK<sub>a</sub> predictions) and unmatched predicted pK<sub>a</sub>s (number of extra pK<sub>a</sub> predictions between 2 and 12. This table is ranked by increasing RMSE. A CSV version of this table can be found in *SAMPL6-supplementary-documents.tar.gz*.

Submission ID	RMSE	MAE	ME	R <sup>2</sup>	m	Kendall's Tau	Unmatched exp. pK <sub>a</sub> s	Unmatched pred. pK <sub>a</sub> s [2,12]
xvxzd	0.68 [0.54, 0.81]	0.58 [0.45, 0.71]	0.24 [-0.01, 0.45]	0.94 [0.88, 0.97]	0.92 [0.84, 1.02]	0.82 [0.68, 0.92]	2	4
gyuhx	0.73 [0.55, 0.91]	0.59 [0.44, 0.74]	0.03 [-0.23, 0.28]	0.93 [0.88, 0.96]	0.98 [0.90, 1.08]	0.88 [0.80, 0.94]	0	7
xmyhm	0.79 [0.52, 1.03]	0.56 [0.38, 0.77]	0.13 [-0.14, 0.41]	0.92 [0.85, 0.97]	0.96 [0.86, 1.08]	0.81 [0.68, 0.90]	0	3
nb017	0.94 [0.72, 1.16]	0.77 [0.58, 0.97]	-0.16 [-0.49, 0.16]	0.88 [0.81, 0.94]	0.94 [0.82, 1.08]	0.73 [0.60, 0.84]	0	6
nb007	0.95 [0.73, 1.15]	0.78 [0.60, 0.97]	0.05 [-0.29, 0.37]	0.88 [0.77, 0.95]	0.84 [0.77, 0.92]	0.79 [0.65, 0.89]	0	13
yqkga	1.01 [0.78, 1.23]	0.80 [0.59, 1.03]	-0.17 [-0.51, 0.19]	0.87 [0.78, 0.93]	0.93 [0.77, 1.08]	0.83 [0.72, 0.91]	0	1
nb010	1.03 [0.77, 1.26]	0.81 [0.61, 1.04]	0.24 [-0.11, 0.59]	0.87 [0.77, 0.94]	0.95 [0.83, 1.08]	0.80 [0.67, 0.90]	0	4
8xt50	1.07 [0.78, 1.36]	0.81 [0.58, 1.07]	-0.47 [-0.82, -0.14]	0.91 [0.84, 0.95]	1.08 [0.94, 1.22]	0.80 [0.68, 0.89]	0	0
nb013	1.10 [0.72, 1.47]	0.80 [0.56, 1.09]	-0.15 [-0.55, 0.22]	0.88 [0.78, 0.95]	1.09 [0.90, 1.25]	0.79 [0.64, 0.90]	0	6
nb015	1.27 [0.98, 1.56]	1.04 [0.80, 1.31]	0.13 [-0.32, 0.56]	0.87 [0.80, 0.93]	1.16 [0.94, 1.34]	0.78 [0.66, 0.86]	0	0
p0jba	1.31 [0.69, 1.73]	1.08 [0.43, 1.72]	-0.92 [-1.72, -0.11]	0.91 [0.51, 1.00]	1.18 [0.36, 1.72]	0.80 [0.00, 1.00]	0	0
37xm8	1.41 [0.93, 1.84]	1.01 [0.68, 1.38]	-0.18 [-0.69, 0.32]	0.83 [0.70, 0.93]	1.16 [0.98, 1.33]	0.70 [0.56, 0.83]	1	1
mkhqa	1.60 [1.13, 2.05]	1.24 [0.90, 1.62]	-0.32 [-0.89, 0.21]	0.80 [0.67, 0.91]	1.14 [0.98, 1.34]	0.64 [0.44, 0.79]	0	6
ttjd0	1.64 [1.20, 2.06]	1.30 [0.96, 1.67]	-0.12 [-0.70, 0.45]	0.81 [0.69, 0.91]	1.2 [1.03, 1.40]	0.65 [0.47, 0.80]	0	5
nb001	1.68 [1.05, 2.37]	1.21 [0.84, 1.68]	0.44 [-0.10, 1.03]	0.80 [0.70, 0.90]	1.16 [0.95, 1.42]	0.72 [0.55, 0.85]	0	7
nb002	1.70 [1.08, 2.38]	1.25 [0.89, 1.70]	0.51 [-0.04, 1.10]	0.80 [0.70, 0.90]	1.15 [0.95, 1.42]	0.72 [0.56, 0.84]	0	7
35bdm	1.72 [0.66, 2.34]	1.44 [0.62, 2.26]	-1.01 [-2.18, 0.13]	0.92 [0.46, 1.00]	1.45 [0.73, 2.15]	0.80 [0.00, 1.00]	0	0
ryzue	1.77 [1.42, 2.12]	1.50 [1.17, 1.84]	1.30 [0.86, 1.72]	0.91 [0.86, 0.95]	1.23 [1.06, 1.41]	0.82 [0.71, 0.91]	0	0
2ii2g	1.80 [1.31, 2.24]	1.39 [1.01, 1.82]	-0.74 [-1.29, -0.15]	0.79 [0.65, 0.89]	1.15 [0.96, 1.37]	0.68 [0.59, 0.82]	0	2
mpwiy	1.82 [1.39, 2.23]	1.48 [1.14, 1.88]	0.10 [-0.54, 0.73]	0.82 [0.70, 0.91]	1.29 [1.12, 1.51]	0.66 [0.49, 0.80]	0	5
5byn6	1.89 [1.50, 2.27]	1.59 [1.24, 1.97]	1.32 [0.84, 1.80]	0.91 [0.85, 0.95]	1.28 [1.10, 1.48]	0.83 [0.72, 0.92]	0	0
y75vj	1.90 [1.50, 2.26]	1.58 [1.21, 1.97]	1.04 [0.46, 1.60]	0.89 [0.79, 0.95]	1.34 [1.16, 1.53]	0.75 [0.57, 0.88]	1	0
w4iyd	1.93 [1.53, 2.28]	1.58 [1.20, 1.98]	1.26 [0.72, 1.76]	0.85 [0.74, 0.92]	1.21 [1.00, 1.40]	0.73 [0.57, 0.85]	0	1
np6b4	1.94 [1.21, 2.71]	1.44 [1.04, 1.94]	-0.47 [-1.08, 0.24]	0.71 [0.60, 0.87]	1.08 [0.81, 1.43]	0.75 [0.62, 0.86]	0	8
nb004	2.01 [1.38, 2.63]	1.57 [1.16, 2.04]	0.56 [-0.10, 1.27]	0.82 [0.72, 0.90]	1.35 [1.15, 1.60]	0.71 [0.54, 0.84]	0	5
nb003	2.01 [1.39, 2.64]	1.58 [1.18, 2.04]	0.52 [-0.14, 1.22]	0.82 [0.73, 0.91]	1.36 [1.16, 1.61]	0.71 [0.54, 0.84]	0	5
yc70m	2.03 [1.73, 2.33]	1.80 [1.48, 2.13]	-0.41 [-1.09, 0.31]	0.47 [0.28, 0.64]	0.56 [0.35, 0.83]	0.53 [0.35, 0.68]	0	27
hytjn	2.16 [1.24, 3.06]	1.39 [0.86, 2.04]	0.71 [0.03, 1.48]	0.45 [0.13, 0.78]	0.62 [0.26, 1.00]	0.47 [0.16, 0.73]	1	27
f0gew	2.18 [1.38, 2.95]	1.58 [1.09, 2.16]	-0.73 [-1.42, 0.04]	0.77 [0.67, 0.89]	1.29 [1.01, 1.63]	0.76 [0.63, 0.86]	0	0
q3pfp	2.19 [1.33, 3.09]	1.51 [0.99, 2.13]	0.59 [-0.10, 1.37]	0.44 [0.13, 0.77]	0.66 [0.27, 1.07]	0.50 [0.20, 0.75]	1	22
ds62k	2.22 [1.62, 2.81]	1.78 [1.34, 2.27]	0.78 [0.06, 1.52]	0.82 [0.70, 0.90]	1.41 [1.20, 1.63]	0.72 [0.55, 0.85]	0	4
xikp8	2.35 [1.94, 2.73]	2.06 [1.66, 2.47]	0.77 [-0.02, 1.58]	0.89 [0.80, 0.95]	1.59 [1.40, 1.81]	0.76 [0.59, 0.89]	1	0
nb005	2.38 [1.79, 2.95]	1.91 [1.44, 2.43]	0.31 [-0.49, 1.15]	0.84 [0.74, 0.91]	1.56 [1.34, 1.82]	0.71 [0.54, 0.83]	0	0
5nm4j	2.45 [1.42, 3.34]	1.58 [0.94, 2.34]	0.05 [-0.80, 1.07]	0.19 [0.00, 0.70]	0.40 [-0.06, 0.81]	0.34 [-0.04, 0.67]	4	1
ad5pu	2.54 [1.68, 3.30]	1.83 [1.24, 2.49]	-0.65 [-1.48, 0.25]	0.76 [0.64, 0.88]	1.43 [1.12, 1.78]	0.77 [0.63, 0.88]	0	0
pwn3m	2.60 [1.45, 3.53]	1.54 [0.83, 2.37]	0.79 [-0.06, 1.77]	0.21 [0.00, 0.63]	0.37 [0.01, 0.78]	0.34 [0.04, 0.63]	1	3
nb006	2.98 [2.37, 3.56]	2.53 [2.00, 3.10]	0.42 [-0.60, 1.47]	0.84 [0.74, 0.92]	1.78 [1.55, 2.06]	0.71 [0.54, 0.84]	0	0
0hxtm	3.26 [1.81, 4.39]	1.92 [1.03, 2.98]	1.38 [0.37, 2.56]	0.08 [0.00, 0.48]	0.28 [-0.17, 0.83]	0.29 [-0.04, 0.61]	3	7

**Table S3. Evaluation statistics calculated for all microscopic pK<sub>a</sub> prediction submissions based on Hungarian match for 8 molecules with NMR data.** Methods are represented via their SAMPL6 submission IDs which can be cross-referenced with Table 1 for method details. There are eight error metrics reported: the root-mean-squared error (RMSE), mean absolute error (MAE), mean (signed) error (ME), coefficient of determination ( $R^2$ ), linear regression slope (m), Kendall's Rank Correlation Coefficient ( $\tau$ ), unmatched experimental pK<sub>a</sub>s (number of missing pK<sub>a</sub> predictions) and unmatched predicted pK<sub>a</sub>s (number of extra pK<sub>a</sub> predictions between 2 and 12). This table is ranked by increasing RMSE. A CSV version of this table can be found in *SAMPL6-supplementary-documents.tar.gz*.

Submission ID	RMSE	MAE	ME	R <sup>2</sup>	m	Kendall's Tau	Unmatched exp. pK <sub>a</sub> s	Unmatched pred. pK <sub>a</sub> s [2,12]
nb011	0.47 [0.30, 0.64]	0.33 [0.22, 0.46]	-0.02 [-0.18, 0.14]	0.97 [0.94, 0.99]	1.01 [0.97, 1.06]	0.90 [0.78, 0.96]	0	36
hdlyq	0.62 [0.47, 0.76]	0.47 [0.33, 0.62]	0.13 [-0.09, 0.34]	0.95 [0.92, 0.97]	0.34 [0.92, 1.09]	0.87 [0.79, 0.93]	0	16
epvmk	0.63 [0.43, 0.81]	0.47 [0.32, 0.63]	-0.02 [-0.25, 0.21]	0.95 [0.89, 0.98]	0.21 [0.91, 1.04]	0.81 [0.68, 0.91]	0	37
xnoe0	0.65 [0.47, 0.82]	0.50 [0.36, 0.66]	-0.1 [-0.32, 0.13]	0.95 [0.89, 0.98]	0.13 [0.92, 1.05]	0.82 [0.69, 0.91]	0	36
gdqeg	0.65 [0.41, 0.89]	0.43 [0.27, 0.62]	0.11 [-0.10, 0.35]	0.94 [0.88, 0.98]	0.35 [0.87, 1.02]	0.83 [0.67, 0.95]	0	53
400ia	0.66 [0.44, 0.86]	0.47 [0.31, 0.64]	0.00 [-0.22, 0.24]	0.94 [0.88, 0.98]	0.24 [0.87, 1.05]	0.85 [0.73, 0.94]	0	35
nb008	0.76 [0.48, 1.02]	0.52 [0.34, 0.73]	-0.08 [-0.37, 0.17]	0.93 [0.85, 0.98]	0.17 [0.79, 0.93]	0.84 [0.73, 0.92]	0	35
ccpmw	0.79 [0.62, 0.94]	0.62 [0.46, 0.80]	-0.17 [-0.44, 0.11]	0.92 [0.86, 0.96]	0.11 [0.82, 1.05]	0.80 [0.67, 0.89]	0	7
0xi4b	0.84 [0.58, 1.07]	0.61 [0.42, 0.83]	0.22 [-0.07, 0.51]	0.92 [0.84, 0.97]	0.51 [0.91, 1.09]	0.81 [0.65, 0.92]	0	32
cwyk	0.86 [0.60, 1.10]	0.62 [0.42, 0.84]	0.13 [-0.16, 0.44]	0.90 [0.82, 0.96]	0.44 [0.86, 1.08]	0.81 [0.64, 0.92]	0	35
ftc8w	0.86 [0.51, 1.17]	0.59 [0.39, 0.83]	0.10 [-0.19, 0.41]	0.90 [0.77, 0.97]	0.41 [0.84, 0.98]	0.75 [0.57, 0.88]	0	35
nxaaw	0.89 [0.56, 1.25]	0.61 [0.41, 0.87]	-0.02 [-0.35, 0.28]	0.89 [0.75, 0.97]	0.28 [0.85, 1.00]	0.79 [0.63, 0.91]	0	29
nb016	0.95 [0.71, 1.18]	0.77 [0.57, 0.98]	-0.23 [-0.56, 0.12]	0.89 [0.83, 0.95]	0.12 [0.82, 1.07]	0.75 [0.62, 0.85]	0	3
kxzt	0.96 [0.56, 1.33]	0.64 [0.41, 0.92]	0.00 [-0.32, 0.36]	0.90 [0.76, 0.97]	0.36 [0.96, 1.13]	0.79 [0.63, 0.91]	0	37
eyetm	0.98 [0.69, 1.27]	0.72 [0.50, 0.97]	-0.32 [-0.65, 0.00]	0.91 [0.86, 0.96]	0.00 [0.94, 1.22]	0.78 [0.64, 0.88]	0	7
cm2yq	0.99 [0.44, 1.54]	0.56 [0.31, 0.90]	0.10 [-0.21, 0.50]	0.91 [0.83, 0.98]	0.50 [0.96, 1.25]	0.89 [0.80, 0.96]	0	36
2umai	1.00 [0.46, 1.54]	0.57 [0.33, 0.91]	0.07 [-0.25, 0.46]	0.91 [0.82, 0.98]	0.46 [0.96, 1.26]	0.87 [0.76, 0.95]	0	36
ko8yx	1.01 [0.76, 1.25]	0.78 [0.56, 1.01]	0.35 [0.01, 0.67]	0.91 [0.82, 0.96]	0.67 [0.96, 1.19]	0.78 [0.64, 0.89]	0	26
wuuvc	1.02 [0.51, 1.53]	0.62 [0.38, 0.93]	0.19 [-0.13, 0.58]	0.88 [0.80, 0.96]	0.58 [0.85, 1.19]	0.90 [0.81, 0.96]	0	36
ktpj5	1.02 [0.51, 1.56]	0.61 [0.37, 0.95]	0.17 [-0.16, 0.57]	0.88 [0.80, 0.96]	0.57 [0.87, 1.22]	0.89 [0.80, 0.96]	0	36
z7fhp	1.02 [0.49, 1.55]	0.61 [0.36, 0.94]	0.08 [-0.24, 0.48]	0.90 [0.82, 0.97]	0.48 [0.97, 1.26]	0.88 [0.80, 0.95]	0	28
arcko	1.04 [0.73, 1.32]	0.77 [0.53, 1.02]	0.37 [0.05, 0.72]	0.89 [0.80, 0.94]	0.72 [0.90, 1.14]	0.78 [0.62, 0.90]	0	24
y4wws	1.04 [0.70, 1.33]	0.74 [0.49, 1.00]	-0.31 [-0.66, 0.05]	0.91 [0.85, 0.96]	0.05 [1.02, 1.26]	0.79 [0.68, 0.88]	0	30
wcvnu	1.11 [0.80, 1.39]	0.84 [0.59, 1.11]	0.28 [-0.10, 0.66]	0.89 [0.77, 0.95]	0.66 [0.98, 1.22]	0.73 [0.54, 0.88]	1	27
8toyp	1.13 [0.61, 1.65]	0.70 [0.42, 1.05]	0.13 [-0.25, 0.56]	0.88 [0.81, 0.96]	0.56 [0.98, 1.29]	0.83 [0.72, 0.92]	0	27
qsicn	1.17 [0.30, 1.65]	0.88 [0.23, 1.54]	-0.76 [-1.54, 0.01]	0.91 [0.46, 1.00]	0.01 [0.52, 1.59]	0.80 [0.00, 1.00]	0	2
wexjs	1.30 [0.95, 1.62]	0.98 [0.68, 1.29]	0.27 [-0.17, 0.74]	0.86 [0.74, 0.93]	0.74 [1.00, 1.29]	0.73 [0.55, 0.86]	0	25
v8qph	1.37 [0.92, 1.79]	0.98 [0.66, 1.34]	-0.15 [-0.64, 0.34]	0.84 [0.70, 0.93]	0.34 [0.97, 1.32]	0.70 [0.55, 0.82]	0	6
w420e	1.57 [1.18, 1.94]	1.23 [0.90, 1.58]	0.09 [-0.48, 0.62]	0.85 [0.76, 0.91]	0.62 [1.08, 1.46]	0.72 [0.60, 0.82]	0	19
6tvf8	1.88 [0.87, 2.85]	1.02 [0.54, 1.66]	0.45 [-0.14, 1.18]	0.51 [0.16, 0.87]	1.18 [0.26, 0.89]	0.61 [0.34, 0.82]	0	55
0wfzo	2.89 [1.73, 3.89]	1.88 [1.17, 2.68]	0.76 [-0.15, 1.77]	0.48 [0.21, 0.75]	1.77 [0.60, 1.37]	0.51 [0.30, 0.70]	0	4
t8ewk	3.30 [1.89, 4.39]	1.98 [1.06, 3.00]	1.32 [0.27, 2.49]	0.07 [0.00, 0.45]	2.49 [-0.17, 0.79]	0.28 [-0.03, 0.6]	0	6
z3btx	4.00 [2.30, 5.45]	2.49 [1.47, 3.65]	1.48 [0.26, 2.86]	0.29 [0.04, 0.60]	2.86 [0.31, 1.44]	0.43 [0.19, 0.63]	0	1
758j8	4.52 [2.64, 6.18]	2.95 [1.85, 4.25]	1.85 [0.48, 3.38]	0.24 [0.02, 0.58]	3.38 [0.20, 1.51]	0.34 [0.08, 0.57]	0	2
hgn83	6.38 [4.04, 8.47]	4.11 [2.52, 5.93]	2.13 [0.07, 4.28]	0.08 [0.00, 0.39]	4.28 [-0.18, 1.43]	0.32 [0.07, 0.56]	0	0

**Table S4. Evaluation statistics calculated for all microscopic pK<sub>a</sub> prediction submissions based on microstate pair match for 8 molecules with NMR data.** Methods are represented via their SAMPL6 submission IDs which can be cross-referenced with Table 1 for method details. There are eight error metrics reported: the root-mean-squared error (RMSE), mean absolute error (MAE), mean (signed) error (ME), coefficient of determination ( $R^2$ ), linear regression slope (m), Kendall's Rank Correlation Coefficient ( $\tau$ ), unmatched experimental pK<sub>a</sub>s (number of missing pK<sub>a</sub> predictions) and unmatched predicted pK<sub>a</sub>s (number of extra pK<sub>a</sub> predictions between 2 and 12. This table is ranked by increasing RMSE. A CSV version of this table can be found in *SAMPL6-supplementary-documents.tar.gz*.

Update this table with dominant microstate accuracy

Submission ID	RMSE	MAE	ME	$R^2$	m	Kendall's Tau	Unmatched exp. pK <sub>a</sub> s	Unmatched pred. pK <sub>a</sub> s [2,12]
nb016	0.52 [0.25, 0.71]	0.43 [0.23, 0.65]	-0.09 [-0.45, 0.30]	0.92 [0.05, 0.99]	0.99 [0.14, 1.16]	0.62 [-0.14, 1.00]	0	3
hdlyq	0.68 [0.49, 0.83]	0.60 [0.39, 0.80]	0.38 [0.02, 0.70]	0.86 [0.47, 0.98]	0.91 [0.45, 1.26]	0.78 [0.4, 1.00]	0	16
nb011	0.72 [0.35, 1.07]	0.54 [0.28, 0.86]	0.45 [0.14, 0.83]	0.86 [0.18, 0.98]	0.93 [0.50, 1.21]	0.64 [0.26, 0.95]	0	36
ftc8w	0.75 [0.52, 0.96]	0.68 [0.50, 0.89]	-0.31 [-0.68, 0.16]	0.87 [0.02, 0.99]	1.12 [-0.11, 1.39]	0.56 [-0.10, 1.00]	0	35
6tvf8	0.76 [0.55, 0.95]	0.68 [0.46, 0.90]	-0.63 [-0.89, -0.35]	0.92 [0.78, 0.99]	0.94 [0.69, 1.41]	0.87 [0.6, 1.00]	0	55
t8ewk	0.96 [0.65, 1.19]	0.81 [0.46, 1.13]	-0.77 [-1.12, -0.38]	0.80 [0.53, 0.96]	0.96 [0.76, 2.26]	0.78 [0.31, 1.00]	1	7
v8qph	0.99 [0.40, 1.52]	0.67 [0.29, 1.17]	-0.09 [-0.75, 0.45]	0.68 [0.11, 0.97]	0.96 [-1.26, 1.16]	0.38 [-0.3, 1.00]	0	6
ccpmw	1.07 [0.78, 1.27]	0.95 [0.60, 1.25]	-0.83 [-1.25, -0.37]	0.74 [0.43, 0.99]	0.95 [0.70, 2.32]	0.89 [0.52, 1.00]	1	8
0xi4b	1.15 [0.75, 1.50]	0.98 [0.63, 1.36]	-0.30 [-0.94, 0.44]	0.77 [0.02, 0.98]	1.26 [0.09, 2.10]	0.51 [-0.14, 1.00]	0	33
cywyk	1.17 [0.88, 1.41]	1.06 [0.74, 1.35]	-0.47 [-1.09, 0.24]	0.73 [0.02, 0.98]	1.15 [-0.04, 2.00]	0.56 [-0.08, 1.00]	0	36
eyetm	1.17 [0.77, 1.52]	1.00 [0.61, 1.41]	-0.89 [-1.38, -0.38]	0.67 [0.30, 0.94]	0.93 [0.65, 2.59]	0.72 [0.29, 1.00]	1	8
nb008	1.26 [0.74, 1.71]	1.09 [0.63, 1.57]	0.47 [-0.40, 1.32]	0.79 [0.01, 0.99]	1.21 [-0.59, 1.85]	0.52 [-0.2, 1.00]	0	38
y4wws	1.41 [0.95, 1.80]	1.22 [0.78, 1.66]	-0.71 [-1.44, 0.06]	0.87 [0.05, 0.98]	1.55 [0.41, 2.02]	0.56 [-0.11, 1.00]	0	31
ktpj5	1.46 [0.83, 2.10]	1.15 [0.67, 1.77]	0.94 [0.29, 1.68]	0.77 [0.01, 0.98]	1.28 [-0.26, 1.60]	0.42 [-0.27, 0.95]	0	37
wuuvc	1.47 [0.84, 2.09]	1.18 [0.70, 1.77]	0.99 [0.36, 1.68]	0.78 [0.01, 0.98]	1.27 [-0.24, 1.58]	0.47 [-0.20, 1.00]	0	37
xnoe0	1.54 [1.09, 2.00]	1.39 [1.02, 1.83]	0.91 [0.11, 1.64]	0.82 [0.01, 0.98]	1.47 [-0.30, 1.79]	0.42 [-0.27, 0.95]	0	37
qsicn	1.58 [1.44, 1.70]	1.57 [1.44, 1.70]	-1.57 [-1.7, -1.44]	1.00 [0.00, 1.00]	1.06		0	2
epvmk	1.66 [1.20, 2.15]	1.50 [1.07, 1.96]	1.12 [0.31, 1.82]	0.82 [0.02, 0.98]	1.47 [-0.21, 1.8]	0.42 [-0.25, 0.95]	0	37
400ia	1.73 [1.33, 2.17]	1.62 [1.29, 2.02]	1.31 [0.53, 1.93]	0.87 [0.03, 0.99]	1.50 [0.07, 1.84]	0.56 [-0.07, 1.00]	0	36
ko8yx	1.75 [1.08, 2.45]	1.44 [0.87, 2.12]	1.38 [0.74, 2.10]	0.97 [0.88, 1.00]	1.66 [1.46, 2.28]	0.91 [0.69, 1.00]	0	27
2umai	1.76 [1.21, 2.35]	1.54 [1.04, 2.11]	1.31 [0.55, 2.03]	0.82 [0.02, 0.98]	1.43 [-0.02, 1.77]	0.47 [-0.17, 0.95]	0	37
cm2yq	1.77 [1.22, 2.36]	1.55 [1.06, 2.12]	1.33 [0.57, 2.04]	0.82 [0.02, 0.98]	1.43 [-0.02, 1.76]	0.47 [-0.17, 0.95]	0	37
nxaaw	1.80 [0.84, 2.80]	1.34 [0.80, 2.18]	0.16 [-0.77, 1.41]	0.59 [0.02, 0.97]	1.37 [-0.08, 2.92]	0.6 [-0.05, 1.00]	0	30
wcvnu	1.90 [1.14, 2.64]	1.57 [0.97, 2.27]	1.44 [0.70, 2.24]	0.97 [0.91, 1.00]	1.78 [1.58, 2.48]	0.91 [0.69, 1.00]	0	27
kxzt	2.00 [1.13, 2.73]	1.64 [1.00, 2.39]	1.64 [1.00, 2.39]	0.83 [0.01, 0.98]	1.42 [-0.21, 1.99]	0.56 [-0.10, 1.00]	0	38
wexjs	2.05 [1.18, 2.93]	1.66 [1.01, 2.47]	1.48 [0.63, 2.39]	0.96 [0.55, 0.99]	1.87 [1.54, 2.29]	0.73 [0.20, 1.00]	0	26
z7fhp	2.14 [1.38, 2.87]	1.80 [1.12, 2.58]	1.28 [0.18, 2.34]	0.78 [0.02, 0.98]	1.71 [-0.41, 2.13]	0.42 [-0.25, 0.95]	0	30
gdqeg	2.38 [1.97, 2.71]	2.25 [1.74, 2.68]	-1.61 [-2.46, -0.37]	0.10 [0.00, 0.98]	0.31 [-0.60, 1.63]	0.29 [-0.45, 1.00]	0	53
8toyp	2.63 [1.89, 3.29]	2.34 [1.59, 3.07]	1.78 [0.47, 2.89]	0.82 [0.02, 0.98]	1.94 [-0.06, 2.39]	0.47 [-0.17, 0.95]	0	29
w420e	2.63 [1.81, 3.53]	2.34 [1.67, 3.18]	1.74 [0.46, 2.92]	0.98 [0.55, 1.00]	2.28 [1.52, 2.41]	0.73 [0.20, 1.00]	0	20
arcko	2.64 [1.23, 3.78]	2.08 [1.10, 3.24]	1.71 [0.44, 3.10]	0.57 [0.04, 0.95]	1.42 [0.56, 2.93]	0.56 [-0.06, 1.00]	0	28
0wfzo	18.72 [11.21, 25.03]	15.80 [9.9, 22.35]	15.09 [8.28, 22.12]	0.09 [0.01, 0.73]	2.35 [-10.18, 8.12]	0.02 [-0.65, 0.66]	0	12
z3btv	22.60 [15.03, 29.00]	19.70 [12.97, 26.69]	19.70 [12.97, 26.69]	0.09 [0.01, 0.72]	2.35 [-10.00, 8.28]	0.02 [-0.66, 0.66]	0	7
758j8	23.76 [16.33, 30.24]	21.00 [14.26, 28.00]	21.00 [14.26, 28.00]	0.09 [0.01, 0.71]	2.35 [-10.34, 8.12]	0.02 [-0.65, 0.65]	0	8
hgn83	27.91 [20.54, 34.52]	25.60 [18.9, 32.64]	25.60 [18.9, 32.64]	0.09 [0.01, 0.72]	2.35 [-10.21, 8.00]	0.02 [-0.65, 0.65]	0	5