

# A fast and interpretable deep learning approach for accurate electrostatics-driven pKa predictions in proteins

Pedro Reis (✉ [pdreis@ciencias.ulisboa.pt](mailto:pdreis@ciencias.ulisboa.pt))

BioISI - Biosystems & Integrative Sciences Institute <https://orcid.org/0000-0003-3563-6239>

Marco Bertolini

Bayer AG

Floriane Montanari

Bayer AG

Walter Rocchia

Istituto Italiano di Tecnologia <https://orcid.org/0000-0003-2480-7151>

Miguel Machuqueiro

BioISI - Biosystems & Integrative Sciences Institute <https://orcid.org/0000-0001-6923-8744>

Djork-Arné Clevert

Bayer AG

---

## Article

**Keywords:** pKa shifts, deep learning, accuracy

**Posted Date:** March 23rd, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-949180/v2>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# 1 A fast and interpretable deep learning approach for 2 accurate electrostatics-driven $pK_a$ predictions in 3 proteins

4 **Pedro B.P.S. Reis**<sup>1,2,3,\*</sup>, **Marco Bertolini**<sup>1</sup>, **Floriane Montanari**<sup>1</sup>, **Walter Rocchia**<sup>3</sup>, **Miguel**  
5 **Machuqueiro**<sup>2,\*</sup>, and **Djork-Arné Clevert**<sup>1,\*</sup>

6 <sup>1</sup>Bayer A.G., Machine Learning Research, Berlin, Germany

7 <sup>2</sup>BiolSI - Biosystems & Integrative Sciences Institute, Faculty of Sciences, University of Lisboa, Campo Grande,  
8 1749-016 Lisboa, Portugal

9 <sup>3</sup>CONCEPT Lab, Istituto Italiano di Tecnologia, Via E. Melen 83, 16152 - Genova, Italy

10 \*e-mail:

## 11 ABSTRACT

Existing computational methods to estimate  $pK_a$  values in proteins rely on theoretical approximations and lengthy computations. In this work, we use a data set of 6 million theoretically determined  $pK_a$  shifts to train deep learning models that are shown to rival the physics-based predictors. These neural networks managed to assign proper electrostatic charges to chemical groups, and learned the importance of solvent exposure and close interactions, including hydrogen bonds. Although trained only using theoretical data, our pKAI+ model displays the best accuracy on a test set of  $\sim 750$  experimental values. Inference times allow speedups of more than 1000 times faster than physics-based methods. By combining speed, accuracy and a reasonable understanding of the underlying physics, our models provide a game-changing solution for fast estimations of macroscopic  $pK_a$  from ensembles of microscopic values as well as for many downstream applications such as molecular docking and constant-pH molecular dynamics simulations.

## 13 Main

14 Many biological processes are triggered by changes in the ionization state of key amino acid side-chains<sup>1,2</sup>.  
15 Experimentally, the titration behavior of a molecule can be measured using potentiometry or by tracking  
16 free energy changes across a pH range. For individual sites, titration curves can be derived from infrared

17 or NMR spectroscopy. Detailed microscopic information can be quickly and inexpensively obtained with  
18 computational methods, and several *in silico*  $pK_a$  calculations have become widely used to provide insights  
19 about structural and functional properties of proteins<sup>3-5</sup>.

20 In Poisson–Boltzmann-based (PB) methods, the solvent is implicitly described while proteins are  
21 represented by point charges in a low dielectric medium<sup>3,4,6,7</sup>. These continuum electrostatics (CE)  
22 methods assume that the  $pK_{\text{half}}$  (the proton binding affinity for a chemical group in a given conformation)  
23 is a good estimate for the macroscopic  $pK_a$  value. This assumption holds when the protein structure  
24 is sufficiently representative of the conformational ensembles corresponding to both protonation states.  
25 Experimentally determined structures exhibit conformations at a minimum energy state, which, in turn,  
26 is related to a specific protonation state. However, biomolecular systems can explore different energy  
27 basins, which may exhibit alternative protonation states. Energy minima can be affected by experimental  
28 conditions, such as temperature, ionic strength and pH. Inaccuracies in  $pK_a$  predictions due to limited  
29 conformational rearrangements can be reduced by increasing the protein dielectric constant from its default  
30 value (2-4), which only accounts for electronic polarisation. The dielectric constant can be used as an  
31 empirical parameter mimicking the effect of the response mechanisms to the local electric field that are  
32 not explicitly taken into account in the model<sup>8-12</sup>. A more computationally expensive approach is to  
33 explicitly include protein motion by sampling conformers via Monte Carlo (MC) or molecular dynamics  
34 (MD) simulations and applying conformational averaging<sup>4,13-15</sup>. Finally, by coupling the sampling of  
35 protonation states at given pH and conformations, constant-pH MD methods<sup>16-20</sup> provide greater insight  
36 into pH-dependent processes<sup>21-25</sup>.

37 As larger data sets of experimental  $pK_a$  values have become available, a new class of purely empirical  
38 methods has been developed. These models rely on statistical fits of empirical parameters weighting the  
39 different energetic contributions into simplified functions. PROPKA<sup>5</sup> is arguably the most popular of such  
40 methods<sup>26</sup>, and has been shown to perform competitively even when compared to higher-level theory  
41 methods<sup>6,27</sup>. The empirical methods are much faster than the physics-based ones although at the cost of  
42 providing less microscopic insights, and their predictive power is unknown on mutations and/or proteins  
43 dissimilar to those composing the training set.

44 The accuracy of most predictors is bound to the estimation of the same quantity, the so-called  $\Delta pK_a$ .  
45 This is the free energy of transferring the ionizable residue from the solvent to the protein, compared to its  
46 neutral counterpart. Since  $pK_a$  values for all amino acids in water have been experimentally determined,  
47 the  $pK_a^{\text{solvent}}$  term can be fixed and, in practice, it can also be adjusted to incorporate systematic errors.  
48 The  $\Delta pK_a$  can be regarded as a sum of mostly electrostatic contributions stemming from the residue  
49 microenvironment. Therefore, an accurate prediction of  $pK_a$  values for a given conformation requires a

50 correct description of the residue interactions with the surrounding protein charges and with the solvent.

51 At their core, deep learning (DL) models are complex non-linear empirical functions fitted to best  
52 map input variables to output properties. Considering chemical properties, such as  $pK_a$  values, which are  
53 dictated by molecular configurations, and provided that enough examples are presented, it is possible  
54 to train a model to map this relationship without the need to solve non-linear equations in 3D or to sort  
55 through the massive space of possible states.

56 In this paper, we have developed two DL-based  $pK_a$  predictors: pKAI and pKAI+, for  $pK_{half}$  and  
57 experimental  $pK_a$  values, respectively. These models have been trained on a database with  $\sim 6$  million  
58  $pK_a$  values estimated from  $\sim 50$  thousand structures using a continuum electrostatics method, PypKa<sup>6</sup>.  
59 pKAI+ displays an unrivaled performance at predicting experimental  $pK_a$  values on a  $\sim 750$  members data  
60 set. Also, pKAI exhibits an accuracy comparable to the PB-based predictor used to generate the training  
61 set while being approximately  $10\text{--}1000\times$  faster. By applying explainable artificial intelligence (XAI)  
62 analysis, we show that these simple models are able to implicitly model most of the required energetic  
63 contributions such as Coulomb interactions, desolvation and hydrogen-bonding. Therefore, the presented  
64 models feature the best characteristics of CE-based methods – accuracy and interpretability – with the  
65 speed provided by empirical approaches.

## 66 Results

67 The main goal of pKAI is to mimic the  $pK_a$  predictive ability of PB-based methods with a significant  
68 computational performance improvement. Our training set is comprised of  $pK_a$  values calculated using  
69 PypKa on a large number of proteins taken from the Protein Data Bank<sup>28</sup>. An elaborate data split was  
70 performed to minimize data leakage from the training set to the validation and test sets (see Methods).  
71 pKAI was designed to be a simple and interpretable model using the minimum structural features that still  
72 capture the electrostatic environment surrounding every titratable residue. The model has been trained on  
73  $\Delta pK_a$  values rather than on absolute values. The  $pK_a$  shift is in fact a more appropriate quantity to learn,  
74 less dependent on the chemical peculiarities of individual amino-acids and more sensitive to the local  
75 electrostatic environment. For example, residues that share a common side-chain chemical group (such as  
76 glutamate and aspartate sharing a carboxylic acid) are influenced by the same environment in a similar  
77 way.

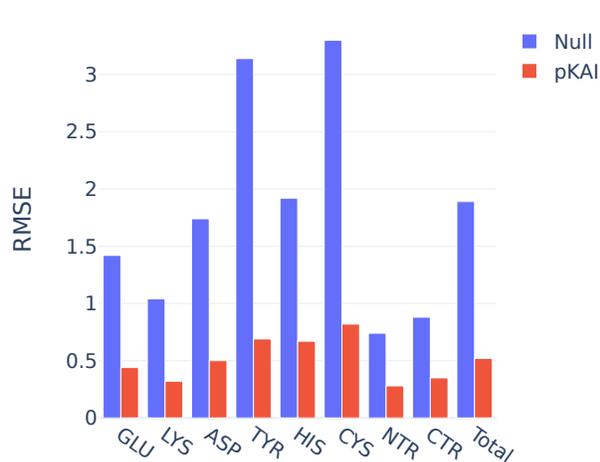
78 We wanted our model to capture the electrostatic dependence between the environment of a residue and  
79 its consequent  $pK_a$  shift while keeping the input layer as small as possible (see Methods). By ignoring all  
80 carbon and hydrogen atoms, we are greatly reducing the dimensionality of our input layer, while retaining

81 most of the information regarding charged particles. There is of course a significant loss of topological  
82 information, although much can be inferred from the positions of the included atoms. In fact, there is  
83 no performance gain when adding solvent exposure measurements (e.g. SASA, residue depth) to the  
84 environment embedding. Considering that solvent exposure entails topological information and that the  
85 model is not able to extract additional information from it, we conclude that it was already estimating, to  
86 some degree, these molecular properties (see Model Explainability subsection).

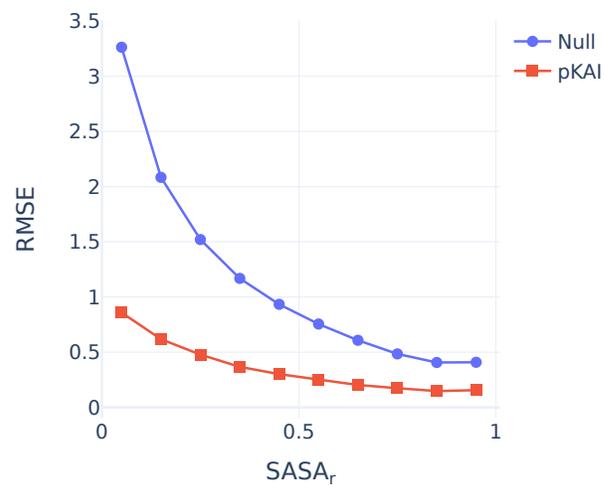
### 87 **pKAI: predicting theoretical $pK_a$**

88 The performance of the model on the test set is reported in Supplementary Table S1 and in Figure 1A. The  
89 null model used for comparison consists of the reference  $pK_a$  value in water for each residue type, and  
90 corresponds to 0 in the  $\Delta pK_a$  scale. Overall, pKAI reproduces the PB-based data with a MAE value of 0.31  
91 and a RMSE of 0.52. However, in this case, we are only predicting theoretical values with a well-defined  
92 relation between structure and  $pK_{\text{half}}$  ( $pK$  value of a single conformation). Experimental  $pK_a$  estimation is  
93 a much more complex task since the  $pK_{\text{half}}$  values corresponding to the different conformations spanned by  
94 the protein should be weighted according to their occurrence probability at equilibrium. The performance  
95 of pKAI is impressive considering the high complexity of the dependence between  $pK_a$  and the site  
96 electrostatic environment, illustrated by the high RMSE value of the Null model (1.89). Some residues are  
97 easier to predict (e.g. LYS and termini residues) while others are more challenging (e.g. CYS and TYR).  
98 This can be explained by their solvent exposure distribution (Figure 1B): well-solvated residues exhibit  
99 small  $\Delta pK_a$  values while more buried ones are more affected by the desolvation effect and establish more  
100 interactions with other residues causing their  $pK_a$  values to shift. There is a clear dependency between the  
101 solvent exposure of a residue, its  $\Delta pK_a$  value and the prediction difficulty (Supplementary Figure S1). The  
102 excellent performance of pKAI is also demonstrated by the fact that most predictions (81.2%) exhibit an  
103 error below 0.5  $pK$  units, which is a sufficient for most use cases.

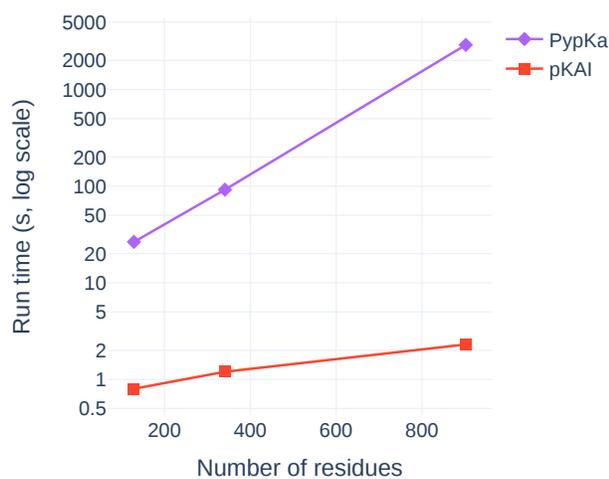
104 The main advantage of DL models is the potential speedup they can provide. Since continuum  
105 electrostatics (CE)  $pK_a$  estimations need to sample thermodynamic equilibrium microstates, several  
106 iterative simulations have to be performed on each protonation state and on the environment of every  
107 residue. On the other hand, pKAI merely needs to apply its learned function over each residue and, as  
108 such, is remarkably faster (Figure 1C). Moreover, the convergence of the CE simulations is harder to  
109 achieve as the protein size increases. Consequently, in PypKa, as the protein size increases, so does the  
110 time required to estimate each  $pK_a$  value. In contrast, the run time of pKAI's DL model has a different  
111 dependence on the protein size. Since the bigger is the protein the larger is the amount of calculations that  
112 can be performed simultaneously, then the less significant becomes the model loading cost and the faster  
113 the average per-residue execution time. This results in sublinear scaling performance and in a speedup



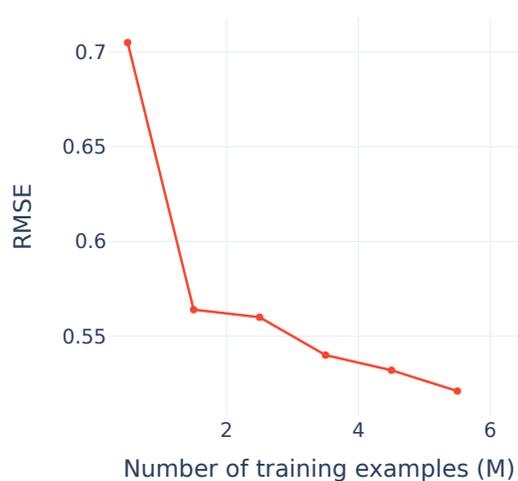
a)



b)



c)



d)

**Figure 1.** A) Comparison between Null model and pKAI RMSE values (values shown in Supplementary Table S1). The Null model is defined as the  $pK_a$  values of the residues in water taken from reference 29. B) Performance at predicting  $pK_{half}$  values dependency on the magnitude of solvent exposure (SASA). The calculations were performed for pKAI and Null model using the PypKa predictions as reference. C) Execution time comparison between PypKa and pKAI (values shown in Supplementary Table S2). This benchmark was executed on a machine with a single Intel Xeon E5-2620 processor. D) Effect of the size of the training set in the model performance on the validation set.

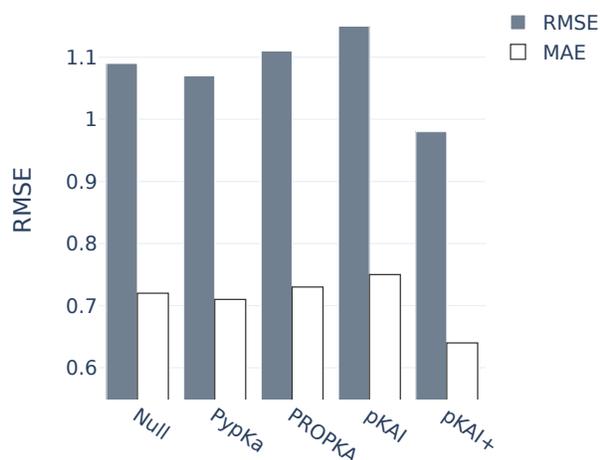
114 over its CE counterpart that can exceed over a thousand times. As such, pKAI is a particularly valuable  
115 tool for dealing with very large systems with thousands of residues where the only added computational  
116 cost stems from the preprocessing of the structure.

117 Another important factor contributing to the high accuracy obtained is the considerable size of the  
118 training set. Despite using the largest repository of experimental protein structures and the largest pK<sub>a</sub>  
119 database available<sup>28</sup>, we show that there is still a correlation between the number of examples in the  
120 training set and the accuracy of the model (Figure 1D). This indicates that our model can still be improved  
121 by providing further examples of pK<sub>a</sub> values. To avoid limiting the scaling rate by the availability of  
122 new experimental protein structures, we can generate new and uncorrelated protein structures using  
123 conformational sampling methods, such as MD and MC. Another advantage of using computational  
124 methodologies is guiding the protein conformational sampling to achieve electrostatic environments that  
125 are underrepresented in the training set. We reserve this development for future work.

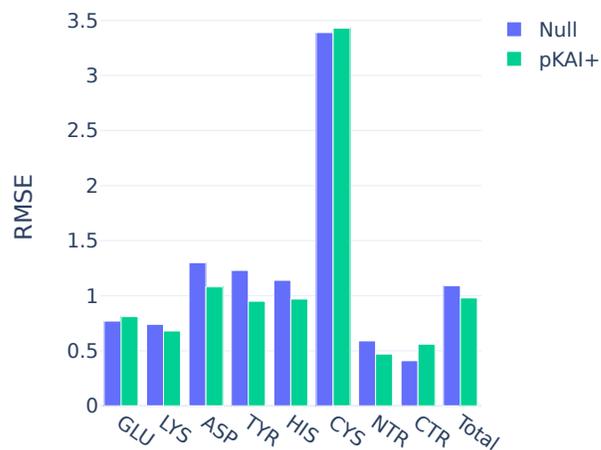
## 126 **pKAI+: Predicting experimental pK<sub>a</sub> values**

127 The main goal of pK<sub>a</sub> predictors, such as PypKa, is to estimate the macroscopic pK<sub>a</sub> value for the  
128 titratable residues using structures (usually experimental ones). Since pKAI aims at reproducing the  
129 pK<sub>half</sub> calculated with PypKa at a fraction of the computational cost, it is not expected to outperform  
130 the PB-based method in predicting experimental values. When using PB to predict experimental pK<sub>a</sub>s,  
131 a higher dielectric constant for the solute is often adopted to compensate for the lack of conformational  
132 flexibility in the method and the lack of representativity of the experimental input structure. A similar  
133 approach can be implemented in pKAI by introducing a regularization weight to the cost function (pKAI+).  
134 This regularization penalizes the magnitude of the  $\Delta pK_a$  prediction. In practice, this procedure biases our  
135 estimates towards the pK<sub>a</sub> values in water, similarly to what is done by the increased solute dielectric  
136 constant in PB-based approaches. It should be noted that pKAI+ has not been trained on experimental  
137 pK<sub>a</sub>, but rather on the same training set as pKAI.

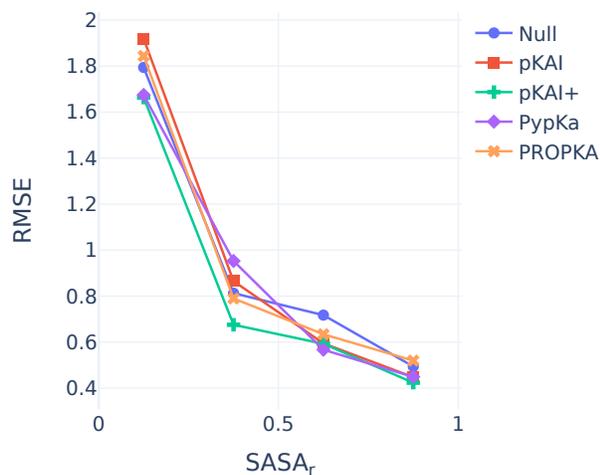
138 To evaluate the performance of our model, we have benchmarked it using a data set of 736 titratable  
139 residues in 97 proteins with experimentally determined pK<sub>a</sub> values (Figure 2A). Remarkably, pKAI+  
140 (RMSE of 0.98) is able to outperform both PypKa (RMSE of 1.07) and PROPKA (state-of-the art empirical  
141 pK<sub>a</sub> predictor, RMSE of 1.11). Furthermore, the improvement over the other methods is significant for  
142 most residue types (Figure 2B), and can be quantified using metrics that are more (RMSE, 0.9 quantile) or  
143 less (MAE, error percentage under 0.5) sensitive to the presence of outliers (Supplementary Table S3).  
144 Cysteine residues are particularly difficult to predict because they naturally occur less frequently and are  
145 more buried than all other titratable residues. This leads to an under-representation of these residues in the



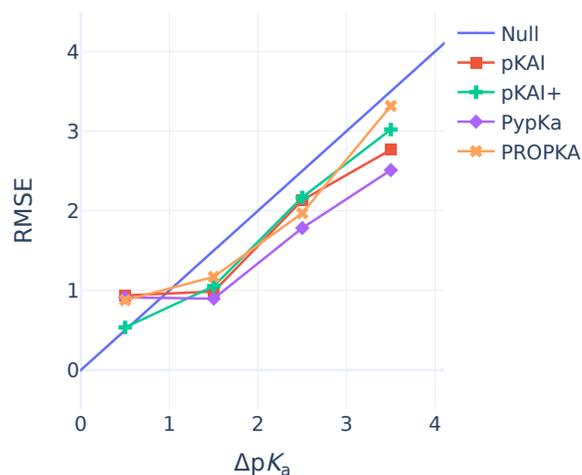
a)



b)



c)



d)

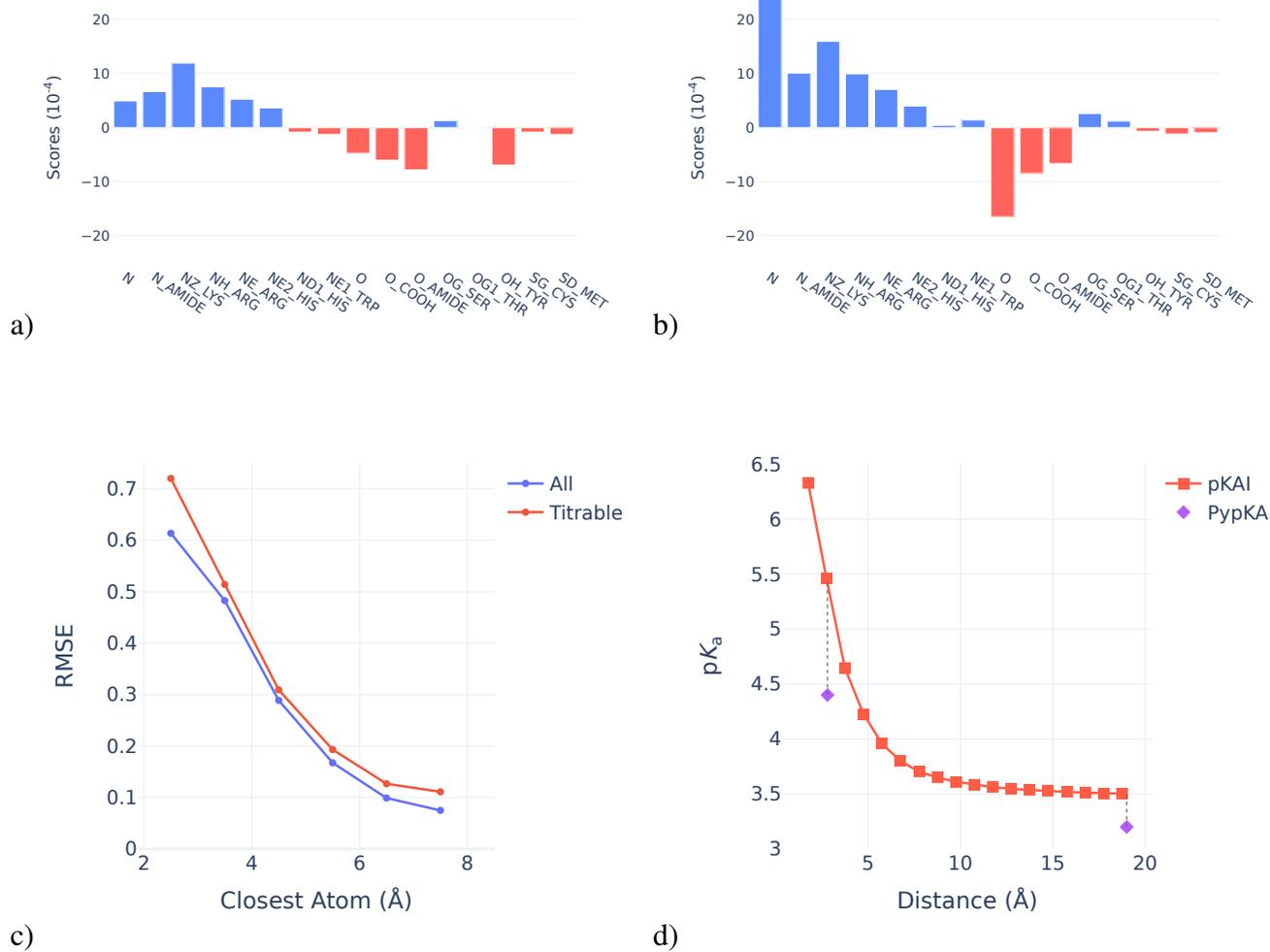
**Figure 2.** A) Experimental pK<sub>a</sub> benchmark of several methods on a data set of 736 residues from 97 proteins (values shown in Supplementary Table S4). The null model values are the pK<sub>a</sub> values of each amino acid substituted in an alanine pentapeptide (Ace-AA-X-AA-NH<sub>2</sub>)<sup>29,30</sup>. B) Comparison between Null model and pKAI+ performance by residue type. C) pKAI+ performance at predicting experimental pK<sub>a</sub> values dependency on the magnitude of solvent exposure (SASA) of the residues. D) Prediction errors of the different models given the experimental pK<sub>a</sub> shift (ΔpK<sub>a</sub>)

146 training set while exhibiting the largest  $pK_a$  shifts. To illustrate the difficulty of this data set, note that  
147 some methodologies are not able to improve on the null model (RMSE of 1.09). The reported deviations  
148 are specific to this data set. Even though our benchmark is one of the largest ever used to validate a  $pK_a$   
149 predictor, it is likely still insufficient to quantify the true accuracy of these methods. Furthermore, besides  
150 being limited, these test sets used for validating new  $pK_a$  predictors tend to always be different. This  
151 makes it very hard to compare methods without rerunning them. In this benchmark, PypKa represents the  
152 PB-based methods like DelPhiPKa<sup>7</sup> or H++<sup>3</sup>. More computationally expensive methods such as MCCE<sup>31</sup>  
153 or constant-pH MD are not represented here. These methods are expected to outperform PB-based  
154 methods, which rely on a single structure, although the exact improvement on this test set is hard to  
155 predict. DeepKa is a recently published convolutional neural network trained on theoretical  $pK_a$  values  
156 from constant-pH MD (CpHMD) simulations<sup>32</sup>. As expected, the CpHMD implemented in the Amber  
157 suite<sup>33</sup> (RMSE of 1.02) outperformed PROPKA (RMSE of 1.12) in their test set which only includes the 4  
158 residues (Asp, Glu, His and Lys) predicted by DeepKa (RMSE of 1.05).

159 Our test set can be divided by solvent exposure (SASA) of the titrating residue.  $pKAI+$  shows  
160 comparable RMSE values to PypKa for both the most solvent exposed and buried residues (Figure 2C).  
161 Interestingly, it is also able to surpass the PB-based model for partially exposed residues. Notably,  
162  $pKAI+$  only improves the PypKa predictions for  $pK_a$  shifts smaller than 1  $pK$  unit (Figure 2D). This  
163 indicates that  $pKAI+$  corrects the  $pK_a$  values of partially exposed residues which are establishing non-  
164 representative interactions in the experimental structure. Since there is a large number of residues with  
165 these characteristics in the test set<sup>28</sup>, the overall performance improvement is significant (Supplementary  
166 Table S4).

## 167 **Model Explainability**

168 The main driving force for  $pK_a$  shifts in proteins is electrostatic in nature. In our model, each atom of  
169 the environment represents the contribution of a chemical group or part of a residue. This individual  
170 contribution towards the final  $\Delta pK_a$  prediction can be estimated (see XAI in the Methods section for further  
171 details) and it is shown in Figure 3A. Remarkably, although our model has been given no information about  
172 atomic charges, it assigns contributions that are in agreement with the expected overall charge of the atom  
173 class. Cationic amine groups (NZ\_LYS; NH\_ARG; NE\_ARG; NE2\_HIS) are clearly assigned positive  
174 scores (i.e. destabilize the protonation of the titratable residue) and are easily distinguishable from the  
175 anionic carbonyl groups (O\_COOH from Glu, Asp and C-termini residues). These scores provide a general  
176 insight into the network's interpretation of each atom and should not be used for more quantitative analysis.  
177 Since the atom score is an averaged measure across the test set, an imbalance of closely interacting atoms  
178 of a specific class can dramatically skew its median contribution.

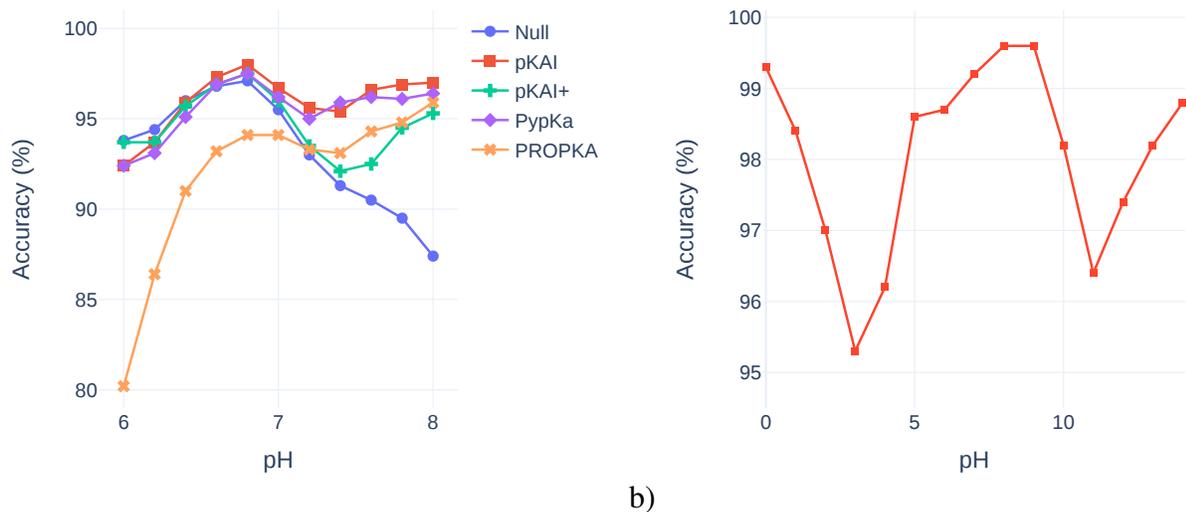


**Figure 3.** Charge scores attributed by pKAI to all considered input atoms classes (Supplementary Table S5) of all atoms (A) and atoms closer than 6  $\text{\AA}$  (B). C) Closest atom influence on pKAI performance. D) Impact of changing the distance of the closest atom on pKAI predictions of residue TYR-315 from structure 2BJU. For reference, we have included PypKa predictions of the same residue in the state presented in the experimental structure (closest distance 2.8  $\text{\AA}$ ) and in a modified structure in which the closest atom is absent.

179 Hydrogen bonds are one of the strongest interactions found in proteins, and, as such, their proper  
180 description is crucial to obtain accurate  $pK_a$  predictions. By comparing Figures 3 A and B we can  
181 observe marked differences between the atom scores at close proximity and those farther away from the  
182 titrating residue. For example, the average score of the very abundant classes of primary amines (N;  
183 N\_AMIDE) and carbonyl groups (O; O\_AMIDE) is greatly diminished when compared to their short-range  
184 contributions, where these become hydrogen donors and acceptors, respectively. The anionic Tyr residue  
185 is perceived to have an overall negative contribution, except when it is close to another titratable residue;  
186 in this case there seems to be no preferred state as it can act both as a donor and as an acceptor – like any  
187 titratable residue. On the other hand, the contribution of neutral non-titrating alcohol groups (OG\_SER;  
188 OG1\_THR) is almost exclusively attributed to their potential to form hydrogen bonds at short range.  
189 Overall, the model is able to capture an astonishing amount of details concerning the physics underlying  
190  $pK_a$  shifts.

191 Beyond the general understanding shown before, hydrogen bond contributions are hard to account  
192 for, compared to other interactions. As shown in Figure 3C the closer another residue (blue curve) is to  
193 the titrating one, the harder for the model is to correctly describe their interaction. The difficulty of the  
194 prediction increases dramatically at the typical distance of hydrogen bonds (2.5-3.2Å). This is even more  
195 marked if one considers interactions established between two titratable residues (red curve). In this case,  
196 the network has to solve for the  $pK_a$  of both residues simultaneously, and in many instances it is unable  
197 to do so. Hence, predicting the contribution of the remaining environment is easier than that of a single  
198 hydrogen bond. This is illustrated in Figure 3D, in which the agreement with the physics-based method is  
199 much higher when the closest atom is removed from the structure rather than when it is kept in its original  
200 position. Although many other profiles can be observed (Supplementary Figure S2), this trend is generally  
201 conserved. Considering that the model did not receive explicit information about hydrogen bonds, it is  
202 quite remarkable that it was able to correlate this type of interaction with larger  $pK_a$  shifts.

203 Solvent exposure is another property that is usually a key contributor to  $pK_a$  shifts. The models are  
204 trained without explicit knowledge of the 3D structure of the protein, and deprived of information regarding  
205 carbon atoms. Nevertheless, they seem to learn about the solvent exposure contribution. We compared  
206 the correlations (Pearson correlation coefficient  $r$  and Spearman's rank correlation coefficient  $\rho$ ) between  
207 the calculated SASA and the  $pK_a$  shifts over the entire test data set. Using the known  $\Delta pK_a$ , we obtained  
208  $r_{\Delta pK_a} = -0.68$ ;  $\rho_{\Delta pK_a} = -0.60$ , while using the predicted  $\Delta pK_a$ , we got  $r_{\text{pred}} = -0.66$ ;  $\rho_{\text{pred}} = -0.62$ .  
209 The similarity between these values indicates that the model has learned the correct correlation between  
210 SASA and the  $pK_a$  shift. Additionally, we trained a model where we provided SASA as an additional  
211 input and observed no performance improvement (data not shown).



**Figure 4.** A) Accuracy of several methods at predicting a representative protonation states derived from experimental  $pK_a$  values. Residues at a pH within 1.5 units of the experimental  $pK_a$  are considered not to have a representative protonation state. B) pKAI accuracy at predicting PypKa derived protonation states.

212 Finally, it is worth mentioning that XAI analysis was a driving factor in the development of pKAI. In  
 213 fact, the importance that the model assigns to each atom class (similar to Figure 3) was pivotal to select  
 214 the final set of atom classes aimed at describing the surrounding environment residues.

## 215 Discussion

216 We have introduced pKAI and pKAI+, two deep learning models to predict theoretical and experimental  
 217  $\Delta pK_a$  values, respectively. pKAI offers unprecedented efficiency, exhibiting a remarkable trade-off between  
 218 accuracy and computational speed, its performance rivaling that of CE-based methods, such as PypKa.  
 219 pKAI could be used as a replacement for such methods, especially when dealing with large proteins or  
 220 applications requiring multiple CE calculations, like constant-pH MD simulations<sup>16–20</sup>. Considering the  
 221 latest advances in sequence to structure predictions<sup>34</sup>, faster methods, such as pKAI, will likely be of  
 222 use as exponentially more structures become available. Furthermore, when optimizing new structures  
 223 for binding to specific targets (e.g. design of enzymes and/or antibodies), it is vital to have an accurate  
 224 prediction of the protonation states.

225 While we strive for optimal accuracy, we are aware that many applications will only require a binary  
 226 decision (hence a qualitative prediction of  $pK_a$  shifts would be sufficient). For example, when selecting the

227 most likely protonation state of a protein, one only needs to predict whether each  $pK_a$  is larger or smaller  
228 than the pH value of interest. The most abundant protonation states obtained from pKAI predictions  
229 are in high agreement with those derived from experiments (Supplementary Figure S3) and outperform  
230 those of PROPKA in a wide range of pH values. Interestingly, PROPKA and the Null model seem to  
231 perform particularly well at extreme pH values. Nevertheless, pKAI is the best model at assigning a fixed  
232 protonation state to a protein at biologically relevant pH values (Fig 4A), arguably the most common task  
233  $pK_a$  predictors are used for. pKAI+ is biased to predict  $pK_a$  values between those of pKAI and the Null  
234 model. While this bias has granted the model an edge on experimental  $pK_a$  estimation, in tasks in which  
235 the Null model does not perform well, pKAI+'s ability is also affected. This can be seen in the biological  
236 range at the more basic pH values.

237 Several other applications only require an estimation of the proton binding affinity using a fixed con-  
238 formation. This quantity, termed  $pK_{half}$ , renders a good prediction of the macroscopic  $pK_a$  when averaged  
239 over a representative ensemble of conformations. From  $pK_{half}$  values, the most abundant/representative  
240 protonation states for a particular conformation can be calculated, improving the realism of methods such  
241 as molecular dynamics<sup>16-20</sup> and molecular docking<sup>35</sup>. pKAI is nearly perfect at mimicking representative  
242 protonation states given by PypKa, being particularly effective at physiological pH, achieving an astound-  
243 ing accuracy of 99.4% (Fig 4B). In a conformational ensemble, there are always many representative  
244 protonation states which differ significantly from the one calculated using the macroscopic  $pK_a$  values.  
245 Therefore, coupling  $pK_{half}$  calculations with conformational sampling techniques is very appealing in  
246 theory but difficult in practice, due to their computational cost. By using pKAI instead of PypKa (or any  
247 other PB-based method), one would drastically decrease the computational overhead (up to 1000 $\times$ ).

248 pKAI does not handle all residues with the same performance. Difficult cases are caused by low  
249 representation in the training set, low solvent exposure, and/or close-by residues providing H bond  
250 interactions. These peculiar environments usually present a high  $\Delta pK_a$  which is not handled very well  
251 by the method. One clear way to improve our models would therefore be to introduce more training  
252 examples. Furthermore, the inclusion of more training data with rare environments would definitely  
253 enhance performance. To better handle interactions with neighboring titratable groups, a change of  
254 environment encoding would be needed. One approach to be explored in future work would be to represent  
255 the whole protein as a graph, and use graph neural networks algorithms to learn the  $\Delta pK_a$  values.

256 Although pKAI excels at predicting  $pK_{half}$  values, its performance is modest when estimating experi-  
257 mental  $pK_a$  values. Inspired by the observation that increasing the dielectric constant in PB-based methods  
258 improves their agreement with experimental results, we have introduced a regularization parameter into the  
259 cost function. Similar to the dielectric constant, this regularization weight biases all predictions towards

260 the residue's  $pK_a$  values in water. The new model, pKAI+, outperforms all methods tested in this work,  
261 including PypKa which was used to create the training set. However, this improvement, while significant  
262 for partially exposed residues which would otherwise exhibit overestimated  $pK_a$  shifts, penalizes the  
263 accuracy of more shifted residues.

264 With pKAI and pKAI+, we are introducing the first deep learning-based predictor of  $pK_a$  shifts  
265 in proteins trained on continuum electrostatics data. The unique combination of speed and accuracy  
266 afforded by our models represents a paradigm shift in  $pK_a$  predictions. pKAI paves the way for accurate  
267 estimations of macroscopic  $pK_a$  values from ensemble calculations of  $pK_{half}$  values, overcoming previous  
268 computational limits. By design, the models were trained using a very simplified view of the surroundings  
269 of the titratable group, accounting only for residues within a 15 Å cutoff, and ignoring all carbon and  
270 hydrogen atoms. This design choice allowed for the models to stay small and fast. Explainability methods  
271 confirmed that this input information was enough for the model to capture crucial features such as  
272 electrostatics, solvent exposure, and environment contributions. The models' initial success introduces  
273 several opportunities for further research, including problem encoding, accounting for conformational  
274 flexibility, interactions with other molecule types (i.e. small molecules, nucleic acids, lipids), and adding  
275 further target properties that could be of interest for other applications.

## 276 Online Methods

### 277 Data set

278 To train our DL models, we used a large publicly available data set of estimated  $pK$  values – the pKPDB  
279 database<sup>28</sup>. This data set of  $\sim 3M$   $pK_a$  values was created by running the PypKa tool with default  
280 parameters<sup>6</sup> over all the protein structures deposited on the Protein Data Bank. The target values to be  
281 fitted by our model are theoretical  $pK_{half}$  values estimated with a PB-based method. This implies that  
282 pKAI will inherit the assumptions and limitations of this class of predictors. Our approach contrasts with  
283 the one usually adopted for training empirical predictors, which entails using experimental values to fit  
284 the model's parameters. The main advantage of this novel approach is that we can train models with  
285 significantly more parameters, such as deep learning ones, since there is now a much larger abundance  
286 of training data. As a comparison, in PROPKA3 only 85 experimental values of aspartate and glutamate  
287 residues were used to fit 6 parameters<sup>5</sup>. Recently, traditional ML models have been trained on  $\sim 1k$   
288 experimental  $pK_a$  values<sup>36,37</sup>. However, testing the real world performance of such methods is difficult  
289 as there is a high degree of similarity among available experimental data. Our larger data set translates  
290 into more diversity in terms of protein and residue types and, more importantly, a wider variety of residue

291 environments. It also helps our models to steer away from the undesired overfitting. Furthermore, the  
292 relation between a structure and our target property is deterministic, contrary to that of experimental  $pK_a$   
293 values, which suffers from the lack of entropic information.

294 The ultimate goal of these methods is to accurately predict experimental  $pK_a$  values and thus, we have  
295 assessed the model's performance with  $\sim 750$  experimental  $pK_a$  values taken from the largest compilation  
296 of experimentally determined  $pK_a$  values of protein residues reported in the literature – the PKAD  
297 database<sup>38</sup>. We compare our experimental results with a null model (attributing to each titratable group  
298 the corresponding  $pK_a$  value in water), PypKa (the method used to generate the training set) and PROPKA  
299 with default settings (the empirical method of reference).

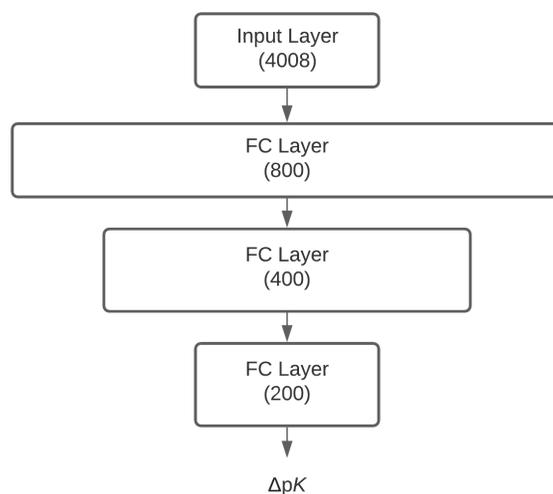
300 Before training our models on our data set, we applied a curated data split (Table 5A) to ensure that  
301 the training, validation, and test sets did not contain proteins with a high degree of similarity and prevent  
302 overfitting. First, we randomly selected 3k proteins from the full data set of  $\sim 120k$  proteins as our holdout  
303 test set of theoretical  $pK_a$  values. The program mmseqs<sup>39</sup> was then used to exclude all proteins containing  
304 at least one chain similar to any of the chains found either in the experimental or in the theoretical test sets.  
305 Chains were considered to be similar if they presented sequence identity over 90%. From the remaining  
306 set of proteins, 3,000 more were randomly assigned to the validation set while the rest became the training  
307 set. Finally, we have excluded similar proteins to those of the validation set from the training set. In the  
308 experimental data set, we have excluded all duplicated proteins, non-exact  $pK_a$  values (e.g.  $> 12.0$ ), and  
309 residues for which PypKa or PROPKA failed to produce an estimate.

## 310 **Model architecture and implementation**

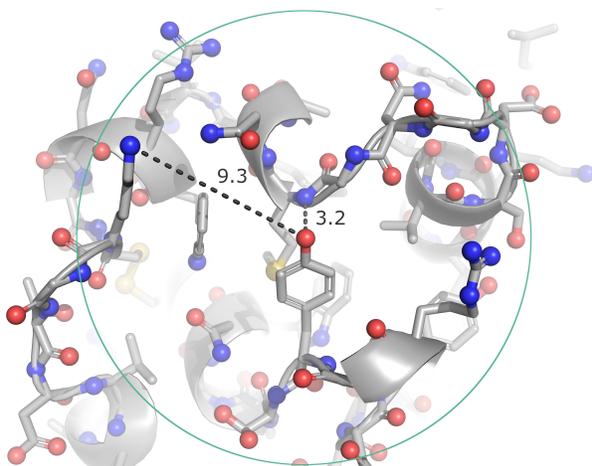
311 pKAI is implemented and trained using PyTorch v1.9.0<sup>40</sup> and PyTorch Lightning v1.2.10<sup>41</sup>. The model  
312 has a simple architecture comprised of 3 fully-connected hidden layers in a pyramidal configuration fitted  
313 to the  $pK_a$  shifts of titratable amino acids (Figure 5B). The simplicity of the architecture is intentional.  
314 pKAI is meant to serve as a proof-of-concept that deep learning models can capture the effect of complex  
315 electrostatic interactions. Describing such interactions would require at least 2 PB calculations per residue  
316 state for the physics-based counterpart (e.g. in PypKa each carboxylic acid has 5 states, hence 10 PB  
317 calculations are required for each Asp/Glu residue).

318 The encoding of the environment of each titratable residue has been simplified to retain only the most  
319 important electrostatic descriptors (Figure 5C). Considering the decay rate of the electrostatic potential, we  
320 decided to truncate the contributions to the environment of a residue by applying a cutoff of 15 Å around  
321 the labile atom(s) of the titratable residue. In practice, this cutoff is slightly smaller for some residue

Split	Proteins	pK values
All Theor.	116.2k	12.6M
Train	56.8k	6.3M
Validation	3.0k	322.4k
Test Theor.	3.0k	325.3k
All Exp.	157	1350
Test Exp.	97	736



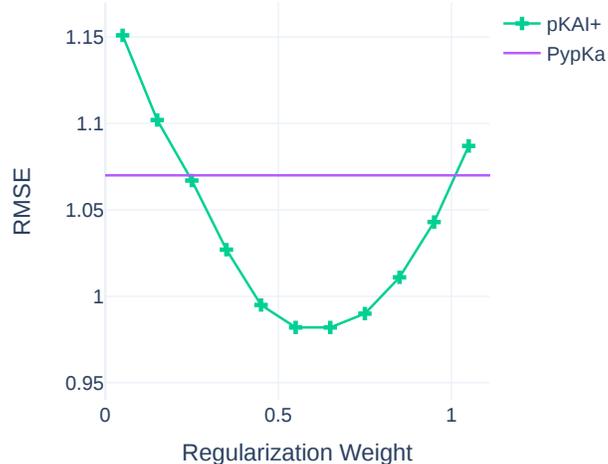
a)



$N = [ (1 / 3.2) 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 ]$   
 $NZ\_LYS = [ 0 0 0 0 0 0 0 0 0 0 0 (1 / 9.3) 0 0 0 0 ]$   
 $TYR = [ 0 0 0 0 0 1 0 0 ]$

c)

b)



d)

**Figure 5.** A) Overview of the data split and similarity exclusion performed on the pKPDB and PKAD databases<sup>28,38</sup>. B) pKAI model architecture. C) Illustration of the titratable amino acid environment encoding. Only Nitrogen, Oxygen and Sulfur atoms (shown as spheres) within a 15 Å cutoff (green circle) are included while all carbon (shown as sticks) and hydrogens (omitted) are ignored. The included atoms are represented by the inverse of their distance to the titratable residue in a OHE vector featuring 16 categories of atom classes (Supplementary Table S5). The titratable residue is represented by a OHE vector comprised of 8 classes. D) pKAI+ regularization weight performance test.

322 environments as the necessary input layer size normalization resulted in the truncation of the closest 250  
323 atoms. A further approximation was made by considering only highly charged atoms and assuming they  
324 are the ones that contribute the most to  $pK_a$  shifts. This simplification can be slightly compensated by using  
325 atom classes instead of charges or element names as they implicitly provide information about adjacent  
326 atoms. The atoms were one-hot encoded (OHE) and, in order to reduce the input layer size, chemically  
327 similar atoms were assigned to the same category (Supplementary Table S5). While carboxylic oxygen  
328 atoms (C-termini OXT; aspartates OD1 and OD2; glutamates OE1 and OE2) and primary amine atoms  
329 (arginines NH1 and NH2) atoms have been merged, others with similar names but different chemical  
330 properties were separated (glutamines OE1 and NE2 from glutamates OE1 and histidines NE2, asparagines  
331 OD1 from aspartates OD1; main chain N from N-termini N).

332 The final 4008-sized input layer consisted of 250 atoms represented by 16 categories OHE classes  
333 concatenated to an 8-dimension OHE vector corresponding to the titrating amino acid. Each atom's OHE  
334 was multiplied by its reciprocal distance to the titrating residues so that this valuable information could be  
335 included without increasing the size of the input layer.

336 pKAI is freely available as a python module that can be installed via pip. The source code can be  
337 found at <https://github.com/bayer-science-for-a-better-life/pKAI>.

## 338 Training

339 Training was performed with mini-batches of 256 examples and the Adam optimizer<sup>42</sup> with a learning  
340 rate of  $1e^{-6}$  and weight decay of  $1e^{-4}$ . Dropout regularization was applied to all fully-connected layers  
341 with the exception of the last one. Hyper-parameter optimization was performed with Optuna<sup>43</sup> using the  
342 performance in the validation set. Training these models takes approximately 10 minutes on an NVIDIA  
343 Tesla M40 24GB, using 16bit precision and an early stopping strategy on the minimization of the cost  
344 function with a delta of  $1e^{-3}$  and patience of 5 steps.

The pKAI model was trained on an MSE cost function while for the pKAI+ we have added a regularization parameter  $\alpha$  to penalize  $\Delta pK_a$  predictions ( $y$ ). Thus, the loss function of pKAI+ becomes

$$J(y_i, \hat{y}_i, \alpha) = (1 - \alpha)(y_i - \hat{y}_i)^2 + \alpha \hat{y}_i^2 \quad (1)$$

345 where  $y_i$  is the true value and  $\hat{y}_i$  the estimation. Different regularization weights were tested on the  
346 validation set to check for overfitting (Figure 5D). While we have selected an  $\alpha$  of 50%, any value in the  
347 40–70% range would lead to a similar improvement.

## 348 XAI Methods

For each input atom feature  $\hat{a} = (a, r_a)$ , where  $a$  indicates the atom class and  $r_a$  the corresponding distance to the liable atom(s) of the titrating residue, we compute the corresponding attribution  $I(\hat{a})$  with the Integrated Gradients (IG) algorithm,<sup>44</sup> as implemented in the `shap` package<sup>45</sup>.  $I(\hat{a})$  measures the sensitivity of the network output with respect to changes in the input  $\hat{a}$ . A large absolute value of  $I(\hat{a})$  indicates that the network assigns high importance to this feature, while the sign of  $I(\hat{a})$  indicates whether the feature contributes positively or negatively to the output. Given that the most important contributions to the  $\Delta pK_a$  are of electrostatic nature, one can try to explain the model inferred charges for each atom class  $a$  by computing the distant-independent score  $C$ :

$$C(a) = \mathbb{E} [r_a^{-1} I_-(\hat{a})] - \mathbb{E} [r_a^{-1} I_+(\hat{a})], \quad (2)$$

349 where  $I_-$  and  $I_+$  are negative and positive  $I$  values, respectively. The  $C$  score of an atom class is thus  
350 the difference between the distance weighted average of examples with negative and positive  $I$  values, over  
351 a large subset (10000 samples) of the test set. The sign of  $C(a)$  in equation 2 resembles the charge that  
352 the network, on average, assigns to a given atom type. For example, if an atom class is being perceived  
353 by the model as contributing negatively to the  $\Delta pK_a$  ( $\mathbb{E} [r_a^{-1} I_-(\hat{a})] > \mathbb{E} [r_a^{-1} I_+(\hat{a})]$  hence  $C(a) > 0$ ), this  
354 would mean that the network learned that this particular atom stabilizes the deprotonated state, which is  
355 characteristic of positively charged groups.

356 The solvent accessible surface area (SASA) values shown in Supplementary Table S1 and in the XAI  
357 subsection have been taken from pKPDB<sup>28</sup>.

## 358 References

- 359 1. Warshel, A. & Åqvist, J. Electrostatic energy and macromolecular function. *Annu. Rev. Biophys.*  
360 *Biophys. Chem.* **20**, 267–298 (1991).
- 361 2. Kim, J., Mao, J. & Gunner, M. Are acidic and basic groups in buried proteins predicted to be ionized?  
362 *J. Mol. Biol.* **348**, 1283–1298 (2005).
- 363 3. Anandakrishnan, R., Aguilar, B. & Onufriev, A. V. H++ 3.0: automating pK prediction and the  
364 preparation of biomolecular structures for atomistic molecular modeling and simulations. *Nucleic*  
365 *Acids Res.* **40**, W537–W541 (2012).
- 366 4. Georgescu, R. E., Alexov, E. G. & Gunner, M. R. Combining conformational flexibility and continuum  
367 electrostatics for calculating pKas in proteins. *Biophys. J.* **83**, 1731 – 1748, DOI: [https://doi.org/10.](https://doi.org/10.1016/S0006-3495(02)73940-4)  
368 [1016/S0006-3495\(02\)73940-4](https://doi.org/10.1016/S0006-3495(02)73940-4) (2002).

- 369 **5.** Olsson, M. H., Søndergaard, C. R., Rostkowski, M. & Jensen, J. H. PROPKA3: consistent treatment  
370 of internal and surface residues in empirical  $pK_a$  predictions. *J. Chem. Theory Comput.* **7**, 525–537  
371 (2011).
- 372 **6.** Reis, P. B. P. S., Vila-Viçosa, D., Rocchia, W. & Machuqueiro, M. PypKa: A flexible python  
373 module for poisson–boltzmann-based  $pK_a$  calculations. *J. Chem. Inf. Model.* **60**, 4442–4448, DOI:  
374 [10.1021/acs.jcim.0c00718](https://doi.org/10.1021/acs.jcim.0c00718) (2020). PMID: 32857502.
- 375 **7.** Wang, L., Zhang, M. & Alexov, E. DelPhiPKa web server: predicting  $pK_a$  of proteins, rnas and dnas.  
376 *Bioinformatics* **32**, 614–615 (2016).
- 377 **8.** Schutz, C. N. & Warshel, A. What are the dielectric “constants” of proteins and how to validate  
378 electrostatic models? *Proteins: Struct. Funct. Genet.* **44**, 400–417 (2001).
- 379 **9.** Voges, D. & Karshikoff, A. A model of a local dielectric constant in proteins. *J. Chem. Phys.* **108**,  
380 2219–2227, DOI: [10.1063/1.475602](https://doi.org/10.1063/1.475602) (1998). <https://doi.org/10.1063/1.475602>.
- 381 **10.** Demchuk, E. & Wade, R. C. Improving the continuum dielectric approach to calculating  $pK_a$ s of  
382 ionizable groups in proteins. *J. Phys. Chem.* **100**, 17373–17387, DOI: [10.1021/jp960111d](https://doi.org/10.1021/jp960111d) (1996).  
383 <https://doi.org/10.1021/jp960111d>.
- 384 **11.** Rocchia, W., Alexov, E. & Honig, B. Extending the applicability of the nonlinear Poisson–Boltzmann  
385 equation: multiple dielectric constants and multivalent ions. *J. Phys. Chem. B* **105**, 6507–6514 (2001).
- 386 **12.** Li, L., Li, C., Zhang, Z. & Alexov, E. On the dielectric “constant” of proteins: Smooth dielectric  
387 function for macromolecular modeling and its implementation in delphi. *J. Chem. Theory Comput.* **9**,  
388 2126–2136, DOI: [10.1021/ct400065j](https://doi.org/10.1021/ct400065j) (2013). PMID: 23585741, <https://doi.org/10.1021/ct400065j>.
- 389 **13.** Beroza, P. & Case, D. A. Including side chain flexibility in continuum electrostatic calculations of  
390 protein titration. *J. Phys. Chem.* **100**, 20156–20163 (1996).
- 391 **14.** Nielsen, J. E. & Vriend, G. Optimizing the hydrogen-bond network in Poisson–Boltzmann equation-  
392 based  $pK_a$  calculations. *Proteins Struct. Funct. Bioinf.* **43**, 403–412, DOI: [10.1002/prot.1053](https://doi.org/10.1002/prot.1053) (2001).  
393 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/prot.1053>.
- 394 **15.** Baptista, A. M. & Soares, C. M. Some theoretical and computational aspects of the inclusion of  
395 proton isomerism in the protonation equilibrium of proteins. *J. Phys. Chem. B* **105**, 293–309 (2001).
- 396 **16.** Baptista, A. M., Teixeira, V. H. & Soares, C. M. Constant-pH molecular dynamics using stochastic  
397 titration. *J. Chem. Phys.* **117**, 4184–4200 (2002).
- 398 **17.** Mongan, J., Case, D. A. & McCammon, J. A. Constant pH molecular dynamics in generalized Born  
399 implicit solvent. *J. Comput. Chem.* **25**, 2038–2048 (2004).

- 400 **18.** Khandogin, J. & Brooks III, C. L. Toward the accurate first-principles prediction of ionization  
401 equilibria in proteins. *Biochemistry-US* **45**, 9363–9373 (2006).
- 402 **19.** Swails, J. M. & Roitberg, A. E. Enhancing conformation and protonation state sampling of hen  
403 egg white lysozyme using pH replica exchange molecular dynamics. *J. Chem. Theory Comput.* **8**,  
404 4393–4404 (2012).
- 405 **20.** Vila-Viçosa, D., Reis, P. B. P. S., Baptista, A. M., Oostenbrink, C. & Machuqueiro, M. A pH replica  
406 exchange scheme in the stochastic titration constant-pH MD method. *J. Chem. Theory Comput.* **15**,  
407 3108–3116, DOI: [10.1021/acs.jctc.9b00030](https://doi.org/10.1021/acs.jctc.9b00030) (2019).
- 408 **21.** Teixeira, V. H., Vila-Viçosa, D., Reis, P. B. & Machuqueiro, M. pK<sub>a</sub> values of titrable amino acids at  
409 the water/membrane interface. *J. Chem. Theory Comput.* **12**, 930–934 (2016).
- 410 **22.** Vila-Viçosa, D., Campos, S. R. R., Baptista, A. M. & Machuqueiro, M. Reversibility of prion  
411 misfolding: insights from constant-pH molecular dynamics simulations. *J. Phys. Chem. B* **116**,  
412 8812–8821 (2012).
- 413 **23.** Morrow, B. H., Koenig, P. H. & Shen, J. K. Atomistic simulations of pH-dependent self-assembly of  
414 micelle and bilayer from fatty acids. *J. Chem. Phys.* **137**, 194902–194902 (2012).
- 415 **24.** Swails, J. M. *et al.* pH-dependent mechanism of nitric oxide release in nitrophorins 2 and 4. *J. Phys.*  
416 *Chem. B* **113**, 1192–1201, DOI: [10.1021/jp806906x](https://doi.org/10.1021/jp806906x) (2009).
- 417 **25.** Reis, P. B., Vila-Viçosa, D., Campos, S. R., Baptista, A. M. & Machuqueiro, M. Role of counterions  
418 in constant-pH molecular dynamics simulations of PAMAM dendrimers. *ACS Omega* **3**, 2001–2009  
419 (2018).
- 420 **26.** Stanton, C. L. & Houk, K. N. Benchmarking pK<sub>a</sub> prediction methods for residues in proteins. *J.*  
421 *Chem. Theory Comput.* **4**, 951–966, DOI: [10.1021/ct8000014](https://doi.org/10.1021/ct8000014) (2008). PMID: 26621236.
- 422 **27.** Lee, A. C. & Crippen, G. M. Predicting pK<sub>a</sub>. *J. Chem. Inf. Model.* **49**, 2013–2033, DOI: [10.1021/](https://doi.org/10.1021/ci900209w)  
423 [ci900209w](https://doi.org/10.1021/ci900209w) (2009). PMID: 19702243.
- 424 **28.** Reis, P. B. P. S., Clevert, D.-A. & Machuqueiro, M. pKPDB: a protein data bank extension database  
425 of pK<sub>a</sub> and pI theoretical values. *Bioinformatics* DOI: [10.1093/bioinformatics/btab518](https://doi.org/10.1093/bioinformatics/btab518) (2021).
- 426 **29.** Thurlkill, R. L., Grimsley, G. R., Scholtz, J. M. & Pace, C. N. pK values of the ionizable groups of  
427 proteins. *Protein Sci.* **15**, 1214–1218 (2006).
- 428 **30.** Grimsley, G. R., Scholtz, J. M. & Pace, C. N. A summary of the measured pK values of the ionizable  
429 groups in folded proteins. *Protein Sci.* **18**, 247–251 (2009).
- 430 **31.** Song, Y., Mao, J. & Gunner, M. Mcce2: improving protein pK<sub>a</sub> calculations with extensive side chain  
431 rotamer sampling. *J. Comput. Chem.* **30**, 2231–2247 (2009).

- 432 **32.** Cai, Z., Luo, F., Wang, Y., Li, E. & Huang, Y. Protein pka prediction with machine learning. *ACS*  
433 *Omega* **6**, 34823–34831, DOI: [10.1021/acsomega.1c05440](https://doi.org/10.1021/acsomega.1c05440) (2021).
- 434 **33.** Huang, Y., Harris, R. C. & Shen, J. Generalized born based continuous constant ph molecular  
435 dynamics in amber: Implementation, benchmarking and analysis. *J. Chem. Inf. Model.* **58**, 1372–1383,  
436 DOI: [10.1021/acs.jcim.8b00227](https://doi.org/10.1021/acs.jcim.8b00227) (2018). PMID: 29949356.
- 437 **34.** Jumper, J. *et al.* Highly accurate protein structure prediction with alphafold. *Nature* **596**, 583–589,  
438 DOI: [10.1038/s41586-021-03819-2](https://doi.org/10.1038/s41586-021-03819-2) (2021).
- 439 **35.** Onufriev, A. V. & Alexov, E. Protonation and pK changes in protein–ligand binding. *Q. Rev. Biophys.*  
440 **46**, 181–209, DOI: [10.1017/S0033583513000024](https://doi.org/10.1017/S0033583513000024) (2013).
- 441 **36.** Gokcan, H. & Isayev, O. Prediction of protein pka with representation learning. *Chem. Sci.* **13**,  
442 2462–2474, DOI: [10.1039/D1SC05610G](https://doi.org/10.1039/D1SC05610G) (2022).
- 443 **37.** Chen, A. Y., Lee, J., Damjanovic, A. & Brooks, B. R. Protein pka prediction by tree-based machine  
444 learning. *J. Chem. Theory Comput.* **0**, null, DOI: [10.1021/acs.jctc.1c01257](https://doi.org/10.1021/acs.jctc.1c01257) (0).
- 445 **38.** Pahari, S., Sun, L. & Alexov, E. PKAD: a database of experimentally measured pK<sub>a</sub> values of  
446 ionizable groups in proteins. *Database* **2019** (2019).
- 447 **39.** Mirdita, M., Steinegger, M. & Söding, J. MMseqs2 desktop and local web server app for fast,  
448 interactive sequence searches. *Bioinformatics* **35**, 2856–2858, DOI: [10.1093/bioinformatics/bty1057](https://doi.org/10.1093/bioinformatics/bty1057)  
449 (2019).
- 450 **40.** Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library. In Wallach, H.  
451 *et al.* (eds.) *Advances in Neural Information Processing Systems 32*, 8024–8035 (Curran Associates,  
452 Inc., 2019).
- 453 **41.** Falcon, W. Pytorch lightning. *GitHub*. Note: <https://github.com/PyTorchLightning/pytorch-lightning>  
454 **3** (2019).
- 455 **42.** Kingma, D. P. & Ba, J. L. Adam: A method for stochastic gradient descent. In *ICLR: International*  
456 *Conference on Learning Representations*, 1–15 (2015).
- 457 **43.** Akiba, T., Sano, S., Yanase, T., Ohta, T. & Koyama, M. Optuna: A next-generation hyperparameter  
458 optimization framework. In *Proceedings of the 25rd ACM SIGKDD International Conference on*  
459 *Knowledge Discovery and Data Mining* (2019).
- 460 **44.** Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks (2017). [1703.01365](https://arxiv.org/abs/1703.01365).
- 461 **45.** Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. In Guyon, I. *et al.*  
462 (eds.) *Advances in Neural Information Processing Systems 30*, 4765–4774 (Curran Associates, Inc.,  
463 2017).

## 464 **Acknowledgements**

465 We would like to thank Paulo Martel and Diogo Vila-Viçosa for the fruitful discussions as well as the  
466 attendees of the Protein Electrostatics meeting ([www.proteinelectrostatics.org](http://www.proteinelectrostatics.org)). We thank Artemi Bendandi  
467 for proofreading the manuscript.

468 PR and MM acknowledge financial support from FCT: SFRH/BD/136226/2018, CEECIND/02300/2017,  
469 UIDB/04046/2020, and UIDP/04046/2020. This work benefited from services and resources provided by  
470 the EGI-ACE project (receiving funding from the European Union's Horizon 2020 research and innovation  
471 programme under grant agreement No. 101017567), with the dedicated support from the CESGA and  
472 IN2P3-IRES resource providers. MB and FM acknowledge funding from the Bayer AG Life Science  
473 Collaboration ("Explainable AI").

## 474 **Author Contributions**

475 PR, FM and MB planned the experiments. PR performed the implementation of pKAI and analysis. MB  
476 performed the explainability experiments. PR, FM and MM conceived and designed pKAI. PR, MM, FM  
477 and MB wrote the paper with the help of DC and WR.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [si.pdf](#)