

1 **Title: Interpretable deep-learning pK_a prediction for small molecule drugs via**
2 **atomic sensitivity analysis**

3

4 **Authors:** Joseph DeCorte¹⁻³, Benjamin Brown^{4,5}, Jens Meiler^{1,2,4-6*}

5 **Affiliations:**

6 ¹Department of Chemical and Physical Biology, Vanderbilt University, Nashville, TN 37232,
7 USA.

8 ²Center for Structural Biology, Vanderbilt University, Nashville, TN 37232, USA.

9 ³Vanderbilt Medical Scientist Training Program Vanderbilt University Medical Center,
10 Vanderbilt University School of Medicine, Nashville TN, 37232-8725, USA.

11 ⁴Department of Chemistry, Vanderbilt University, Nashville, TN 37232-8275, USA.

12 ⁵Center for Applied AI in Protein Dynamics, Vanderbilt University, Nashville, TN 37232-8725,
13 USA

14 ⁶Institute for Drug Discovery, Leipzig University Medical School, Leipzig, SAC 04103,
15 Germany

16 *Corresponding author. Email: jens.meiler@vanderbilt.edu

17 **ORCID**

18 Joseph DeCorte: orcid.org/0000-0002-0131-6176

19 Benjamin Brown: orcid.org/0000-0001-5296-087X

20 Jens Meiler: orcid.org/0000-0001-8945-193X

21

22 **ABSTRACT**

23 Machine learning (ML) models play a crucial role in predicting properties essential to drug
24 development, such as a drug's logscale acid-dissociation constant (pK_a). Despite recent
25 architectural advances, these models often generalize poorly to novel compounds due to a
26 scarcity of ground-truth data. Further, these models lack interpretability, in part due to a
27 dependence on explicit encodings of input molecules' molecular substructures. To this end,
28 atomic-resolution information is accessible in chemical structures by observing model response
29 to atomic perturbations of an input molecule; however, no methods exist that systematically
30 utilize this information for model and molecular analysis. Here, we present BCL-XpKa, a
31 substructure-independent, deep neural network (DNN)-based pK_a predictor that generalizes well
32 to novel small molecules. BCL-XpKa discretizes pK_a prediction from a regression problem into
33 a multitask-classification problem, which accumulates data for prediction at biologically relevant
34 pH values and records the model's uncertainty in its prediction as a discrete distribution for each

35 pK_a prediction. BCL-XpKa outperforms modern ML pK_a predictors and accurately models the
36 effects of common molecular modifications on a molecule's ionizability. We then leverage BCL-
37 XpKa's substructure independence to introduce atomic sensitivity analysis (ASA), which quickly
38 decomposes a molecule's predicted pK_a value into its respective atomic contributions without
39 model retraining. When paired with BCL-XpKa, ASA informs that BCL-XpKa has implicitly
40 learned high-resolution information about molecular substructures. We further demonstrate
41 ASA's utility in structure preparation for protein-ligand docking by identifying ionization sites in
42 97.8% and 83.4% of complex small molecule acids and bases. We then apply ASA with BCL-
43 XpKa to understand the physicochemical liabilities and guide optimization of a recently
44 published KRAS-degrading PROTAC.

45

46 INTRODUCTION

47 Predicting a drug's behavior in the body is a key challenge in computational drug development.
48 For example, accurate prediction of compounds' bioavailability could support early modification
49 or termination of nonviable lead molecules, thereby saving years of time and millions of dollars
50 on research and development. The demand for fast and accurate predictions of a drug's
51 quantitative structure-activity and structure-property relationships (QSAR, QSPR) has
52 skyrocketed as our access to synthesizable chemical space approaches one trillion molecules¹.
53 While advances in machine learning have improved prediction accuracy, the small amount of
54 publicly available, high-quality experimental data for training often leads to overfitting and
55 prevents generalizability²⁻⁴. Further, QSPR model interpretability is often poor despite the
56 relatively more intuitive input (chemical structures) than in other fields of computational biology
57 (e.g., transcriptomics data). As such, additional explorations into architectures that can efficiently
58 train on chemical data, as well as general methods to interpret these models' outputs, are
59 warranted.

60 One of the most critical properties to a drug's downstream efficacy is its ionizability at
61 physiologic pH values, which depends on the drug's logscale acid-dissociation constant (pK_a
62 value)^{4,5}. Quantum mechanical (QM) methods now calculate pK_a with experimental accuracy
63 and are extremely valuable to late-stage drug development. However, small-molecule drug
64 development often begins with virtual high-throughput screening (vHTS) of billions of
65 compounds, and QM methods are too computationally expensive to assist meaningfully in vHTS.
66 As such, scientists have made tremendous investment in ML-based QSAR/QSPR predictors for
67 faster – though potentially less accurate – prediction of physicochemical properties like pK_a.

68 These ML methods generally embed molecules using molecular fingerprints, which are two-
69 dimensional (2D)- or 3D chemical substructures centered around each atom in the molecule.
70 Recently, groups have realized significant gains in prediction accuracy with graph neural
71 networks (GNNs), which embed molecules as a graph in addition to standard chemical
72 descriptors⁶⁻⁸. Improvements in molecular featurization strategies and network architecture have
73 driven state-of-the-art pK_a prediction accuracy to within 0.75-1.00 pK_a units of experimental
74 values.

75 Despite these advances, several limitations persist in ML-based pK_a prediction and QSPR
76 prediction generally. First, all ML-based pK_a predictors to date use regression. While regression
77 is the natural setting for predicting continuous values, small training set sizes restrict the
78 accuracy of regression outputs, particularly at extreme – but still physically relevant – values.
79 Second, many of these models directly encode common molecular substructures in their feature
80 set⁹. This strategy may limit generalizability by preventing complete consideration of each
81 atom's local context in a molecule. For example, amide Nitrogen atoms are generally not
82 ionizable at physiologic pH values, but their acidity can be greatly increased by appending
83 neighboring electron-withdrawing groups (e.g., diacetamide). Therefore, encoding atom-specific
84 local environments may increase generalizability.

85 Finally, model explainability and interpretability in computational chemistry largely focus on
86 feature-set-level analysis, encompassing step-wise^{10,11}, feature-masking¹², feature-set-
87 perturbation¹³, feature-attribution¹⁴, and response-randomization¹⁵ methods. However, feature-
88 set-level analysis is often slow and unintuitive, as feature sets are often large and complex.
89 Recently, ML frameworks have been developed that utilize direct chemical representations (as
90 SMILES strings) to identify important substructures for transformer predictions¹⁶, but no group
91 to date has leveraged perturbations to the input chemical structure itself to gain insights into
92 model learning and output. Indeed, chemical structures are unique in computational biology in
93 that they can be perturbed in consistent, physical meaningful ways. With an appropriately
94 constructed feature set, measuring a model's response to these perturbations would provide
95 granular details into both model learning and molecular hotspots for prediction in real time,
96 without the need for model retraining. For example, replacing a molecule's acidic carboxylic
97 acid functional group with an inert ketone increases the pK_a from ~4 to ~20, thereby
98 demonstrating the carboxylic acid's importance to acidity. As a counterexample, increasing or
99 decreasing the expression of a gene in a transcriptomics-based predictor of cellular activity may
100 not be physically meaningful, as gene expression is highly dependent on the network of
101 expressed genes in a cell/tissue. This presents an underutilized opportunity for computational
102 chemists to gain valuable insights into model learning, performance, and molecular structure.

103 To address these limitations, we present BCL-XpKa, a substructure-independent multitask
104 classifier for rapid and accurate pK_a prediction built in the Biology and Chemistry Library
105 (BCL), an open-source cheminformatics platform developed and maintained by our lab. We use
106 pK_a prediction to illustrate that discretizing continuous problems in chemical biology into
107 multitask classification problems can increase prediction accuracy without meaningful
108 information loss, thereby circumventing many of the problems associated with regression
109 models. We couple BCL-pKa with a novel method of atomic sensitivity analysis (ASA), which
110 provides unprecedented, atomic-level insights into which regions of a molecule are most
111 important for the model's final prediction. We demonstrate ASA's utility in probing model
112 learning, as well as its direct applicability to introducing targeted modifications in molecules that
113 reduce ionizability. Importantly, ASA is relevant to all forms of QSAR/QSPR prediction and can
114 be easily implemented to existing algorithms.

115 BCL-XpKa is a multi-layer perceptron (MLP) that embeds molecules using 2D chemical
116 descriptors features that only encode information about each atom's local environment (up to 1
117 bond away) in the molecule¹⁷. This scheme is substructure independent and enables increased
118 sensitivity to atom-level perturbations, which is particularly important for hit-to-lead and lead
119 optimization in late-stage drug development. Small molecule drugs often have both basic and
120 acidic regions that vary greatly in pK_a values (e.g., amino acids have both an acidic carboxylic
121 acid and a basic free amine group). To account for this, we trained two models: one to predict a
122 molecule's most acidic pK_a value (BCL-XpKaAcid), and one to predict its most basic pK_a value
123 (BCL-XpKaBase). This is a common practice in modern pK_a prediction. We trained BCL-XpKa
124 on datasets of both predicted and experimental pK_a values, and we evaluate our models on an
125 external test set of challenging acids and bases with experimental pK_a values. We find that our
126 model has competitive accuracy and reduced substructure dependence than state-of-the-art pK_a
127 predictors, including GNN-based models.

128 Overall, the work presented here has the following significant contributions to the field:

- 129 - We developed a novel, substructure-independent framework for QSPR prediction that
130 uses local atomic environment embeddings and replaces regression with multitask
131 classification, using pK_a prediction to illustrate competitive performance with modern
132 ML models.
- 133 - We developed a method that rapidly assesses QSPR model learning and provides atomic-
134 level insights to molecular ionizability without requiring model retraining
- 135 - We integrate these two tools in a workflow for lead optimization and apply it to optimize
136 a pan-KRAS degrading Proteolysis Targeting Chimera (PROTAC).

137

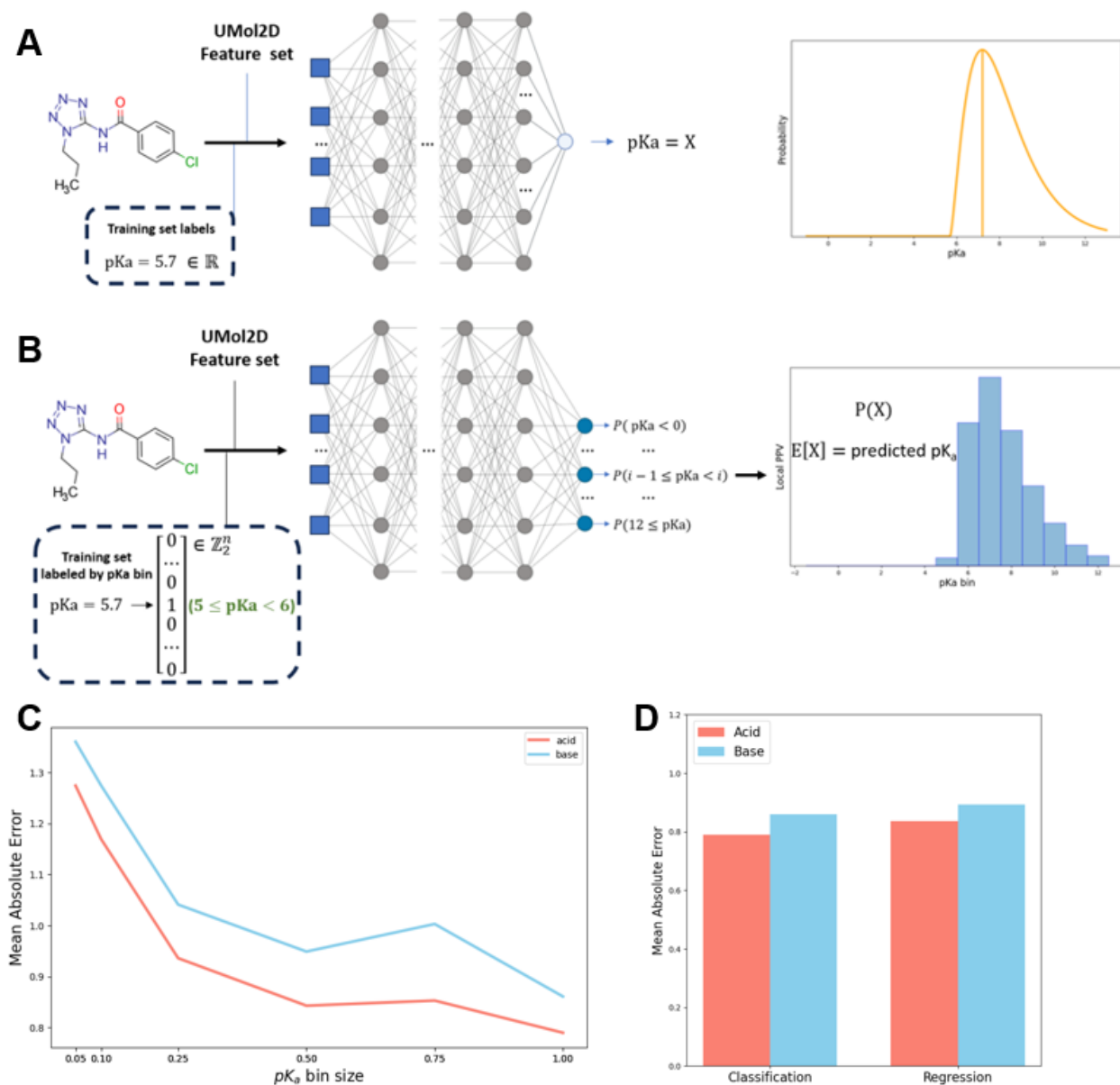
138 RESULTS

139 Multitask Classification is competitive with Regression for small-molecule pK_a prediction

140 Regression models naturally dominate machine-learning (ML) approaches to QSPR prediction.
141 However, regression models require large, high-quality datasets to train properly, and
142 understanding the model's uncertainty in its prediction is challenging. Here, we construct BCL-
143 XpKa, a multilayer perceptron (MLP)-based pK_a prediction tool that transforms a continuous
144 prediction problem into a multitask classification problem. Rather than predict pK_a values on a
145 continuous range, BCL-XpKa predicts the probability that a molecule's pK_a lies within a certain
146 range (Figure 1A-B). The expected value of this probability distribution corresponds to BCL-
147 XpKa's predicted pK_a for a molecule, and the variance of this distribution directly informs its
148 confidence in that prediction.

149 To train BCL-XpKa, pK_a values in the training set were converted into vectors in \mathbb{Z}_2^n , where n is
150 the number of bins the continuous interval has been divided into, and a 1 at position i indicates
151 the molecule's pK_a lies in bin i (Figure 1B). BCL-XpKa models were trained to predict the most
152 acidic and most basic pK_a values of a molecule using training data from ChEMBL augmented
153 with negative data (i.e., nonionizable molecules).

154 Binning training data in this way necessarily leads to some information loss. Using an external
 155 test set, we demonstrate that mean absolute error (MAE) increases as the number of bins
 156 increases, as increasing the number of bins reduces the amount of data in each bin for training
 157 (Figure 1C). A bin size of 1 pK_a unit yielded the lowest MAE on this test set and is used in the
 158 BCL-XpKa production model. For both acids and bases, BCL-XpKa marginally outperforms the
 159 best-performing regression models trained on the same data (0.79 vs 0.83 for acids, 0.86 vs 0.92
 160 for bases, Figure 1D). Additional model details, including an evaluation of model
 161 hyperparameters, can be found in Supplemental Figure 1.



162
 163 **Figure 1 BCL-XpKa model description and internal performance (A-B)** Overview of the
 164 regression and MTC architectures for pK_a prediction. BCL-XpKa utilizes the MTC architecture
 165 with bin size of 1 pK_a unit. (C) MTC model performance by pK_a bin size on acids (red) and

166 bases (blue). (D) Model performance comparison between best performing MTC (BCL-XpKa)
167 and regression architectures on acids (red) and bases (blue).

168

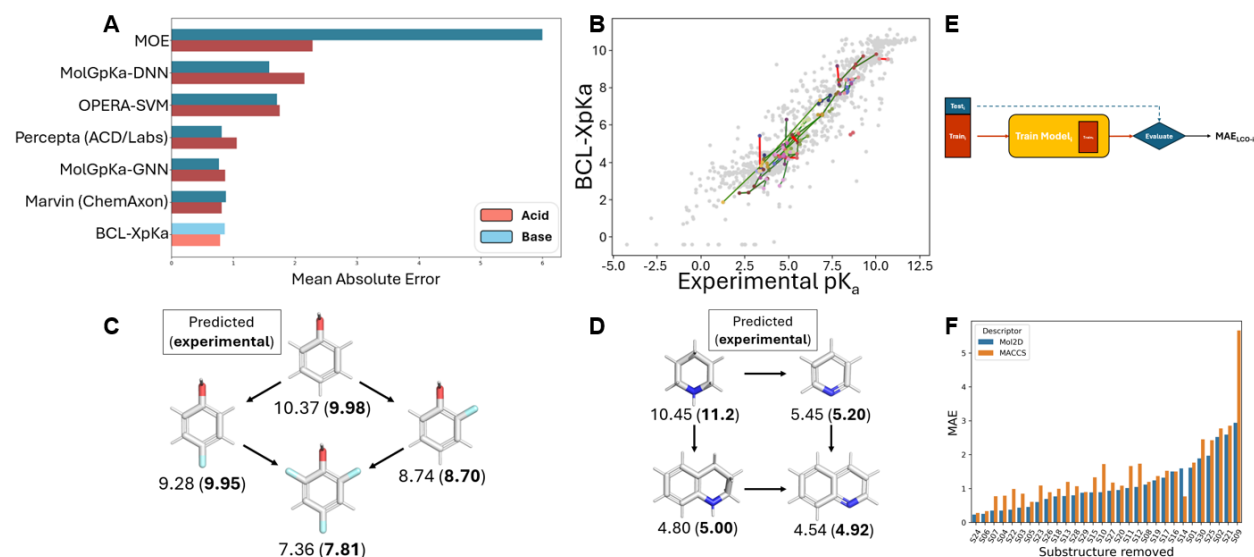
169 **BCL-XpKa accurately captures complex trends in ionizability for druglike small molecules**

170 We then compared BCL-XpKa's performance head-to-head against several state-of-the-art pK_a
171 predictors of varying architecture using a challenging test set external to BCL-XpKa's training
172 data^{7,18-21}. Despite its relatively simple architecture, BCL-XpKa achieves competitive
173 performance to the best machine-learning (MolGpKa, graph-convolutional neural network) and
174 rule-based pK_a predictors on both acids and bases, with BCL-XpKaAcid and BCL-XpKaBase
175 achieving mean absolute error (MAE) of 0.79 and 0.86, respectively (Figure 2A).

176 Beyond accurate pK_a prediction, correctly predicting the effect of small perturbations to a lead
177 molecule's ionizability is of key importance in drug development²². To assess BCL-XpKa's
178 sensitivity to such changes, we identified 71 pairs of molecules in our test set that vary by a
179 slight modification, such as the replacement of an amide with an ester. BCL-XpKa correctly
180 predicts the direction of pK_a change in 81.7% of these pairs (Figure 2B). To illustrate this effect,
181 BCL-XpKa correctly predicts the inductive effect of electron-withdrawing groups on acidity
182 using a series of phenol derivatives (Figure 2C). Here, fluorination at the ortho position increases
183 acidity more than fluorination at the para position (8.74 vs 9.28), and substitution with multiple
184 fluorine atoms has a greater effect than monosubstitution (7.36). While phenol was in our
185 training data, the remaining molecules were not. For bases, BCL-XpKa correctly predicts the
186 complex impact of aromaticity on nitrogen basicity in a series of piperidine derivatives relevant
187 to drug development. Introducing a neighboring phenyl group reduces predicted pK_a from 10.45
188 (true pK_a 11.2) to 4.80 (5.00). Similarly, aromatization of piperidine to pyridine decreases
189 predicted pK_a to 5.45 (5.20), and appending the same phenyl group to produce quinoline reduces
190 pK_a to 4.54 (4.92) (Figure 2D).

191 Finally, substructure independence is critical to QSPR model generalizability to novel
192 compounds. As described above, BCL-XpKa embeds molecules solely using the 1-bond-length
193 neighborhoods of each atom to limit substructure dependence. To investigate this strategy's
194 impact, we subset BCL-XpKa's training set according to 30 ionizable functional groups. We
195 iteratively retrained BCL-XpKa leaving each substructural class out, then tested each model on
196 its withheld substructural class (Figure 2E). BCL-XpKa demonstrates robust performance on this
197 leave-class-out (LCO) test, with an average MAE of 1.1 pK_a units across all LCO models.
198 Training on MACCS descriptors rather than Mol2D yielded systematically worse results and an
199 average MAE of 1.46 pK_a units (Figure 2F).

200



201

202 **Figure 2: BCL-XpKa external performance and molecular series** (A) Performance of various
 203 pK_a predictors on an external test set of acids (red) and bases (blue). (B) BCL-XpKa prediction
 204 vs experimental pK_a value for families of related druglike molecules. Green denotes correct
 205 change in predicted pK_a due to chemical modification, and red denotes incorrect change in
 206 predicted pK_a. (C-D) Example molecular families from (B). Predicted and experimental pK_a
 207 values provided. (E) Schematic for LCO testing. Test_i and Train_i denote the subsets of the
 208 training set that contain or do not contain, respectively, substructure i. Model_i was trained with
 209 Train_i and evaluated on Test_i to give MAE_{LCO-i}. (F) LCO performance of BCL-XpKa (blue) vs an
 210 equivalent model trained with a MACCS-based descriptor set (orange).

211

212 Atomic sensitivity analysis provides actionable, atomic-resolution information on model 213 predictions

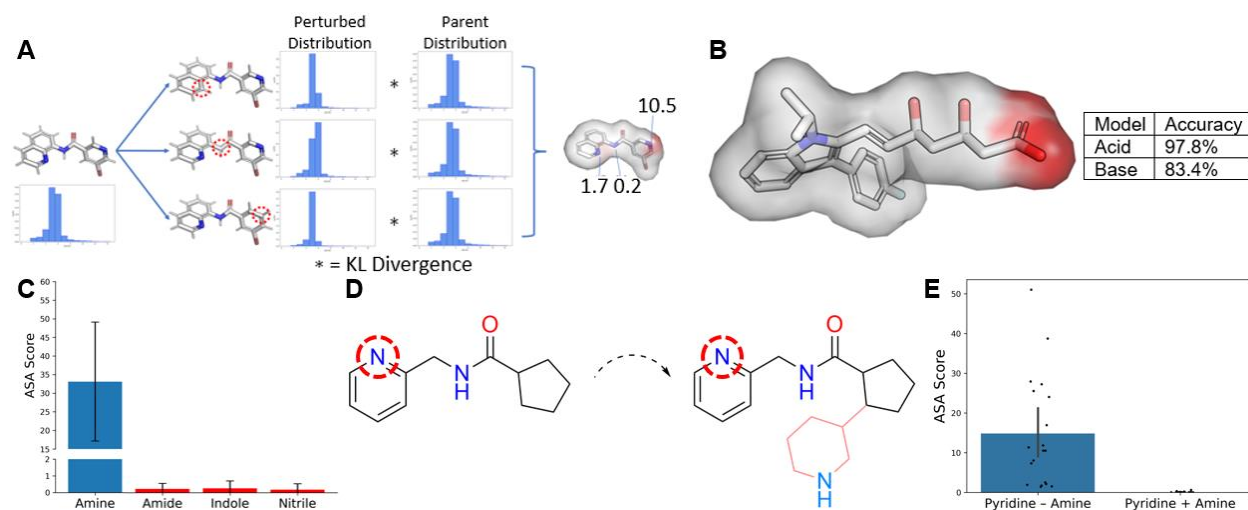
214 Computational chemistry currently lacks rapid, reliable tools for interpreting ML model
 215 predictions²³. Such atom-, substructure-, or pharmacophore-level information could accelerate
 216 computer-aided drug development on multiple fronts, from assisting in model training, to
 217 preparing and filtering molecules in virtual high-throughput screens, to guiding lead compound
 218 modification in lead optimization. To address this deficit, we developed atomic sensitivity
 219 analysis (ASA, Figure 3A), and we demonstrate its utility in decomposing BCL-XpKa's
 220 predictions to assess model learning, identify ionization sites in complex small molecules, and
 221 guide lead optimization efforts by reducing molecular ionizability.

222 ASA compares an ML model's prediction on a parent molecule before and after some
 223 perturbation. Here, we sequentially replace heteroatoms in the parent molecule with correctly
 224 hybridized carbons, then generate probability distributions for each with BCL-XpKa. The final
 225 ASA score is a scaled version of the KL-Divergence between these parent and perturbed pK_a
 226 distributions (see *Methods*). Importantly, we anticipate that, with careful feature-set selection,

227 this replace-rescore-compare scheme will generalize well to substructure- and pharmacophore-
228 level perturbations (Figure 3A).

229 We benchmarked ASA on BCL-XpKa on all test-set molecules with nontrivial ionization sites.
230 We first hypothesized that perturbing an acid's most acidic hydroxyl group or a base's most
231 basic Nitrogen atom would have the most significant impact on the predicted pKa, and therefore
232 the largest ASA scores. We tested this hypothesis by performing ASA on the Oxygen atoms in
233 the acid set and the Nitrogen atoms in the base set and considering only the atom with the
234 maximum ASA score in each molecule. This strategy correctly identifies 97.8% of the most
235 acidic Oxygen atoms in the acid set and 83.4% of the most basic Nitrogen atoms in the base set,
236 thereby demonstrating ASA's potential utility in high-throughput structure preparation (Figure
237 3B).

238 This benchmark revealed surprisingly consistent ASA scores for each atom in the substructures
239 that recurred throughout the test set. For example, free amines are the most basic group in 33.9%
240 of our experimentally characterized bases, and in each of these molecules the amine Nitrogen
241 atom dominates the molecule's ASA scores (33.1 +/- 16.0). These scores were significantly
242 higher than average scores for Nitrogen atoms in amide (0.225 +/- 0.332), indole (0.261 +/-
243 0.444), and nitrile groups (0.180 +/- 0.357) ($p < 0.001$), functional groups which are not
244 traditionally ionizable at physiologic pH and which dominated 0% of the test-set ASA scores
245 (Figure 3C). Further, molecules where the dominant ionizable group is less basic than typical
246 amines also demonstrated consistent ASA scores, and this ASA dominance was reliably ablated
247 by the introduction of an amine functional group (Figure 3D-E).



248
249 **Figure 3: Atomic sensitivity for molecular analysis** (A) Schematic of the ASA protocol. Parent
250 and Perturbed distributions refer to the localPPV distributions output by BCL-XpKa. (B) ASA
251 accuracy at detecting the most acidic Oxygen atom and most basic Nitrogen atom in all non-
252 trivial test-set molecules. (C) ASA scores of positive- and negative-control substructures for
253 BCL-XpKaBase decomposition. Blue denotes the positive control, red denotes the negative
254 controls. (D-E) Modulation of pyridine Nitrogen ASA score by addition of an amine group.

255

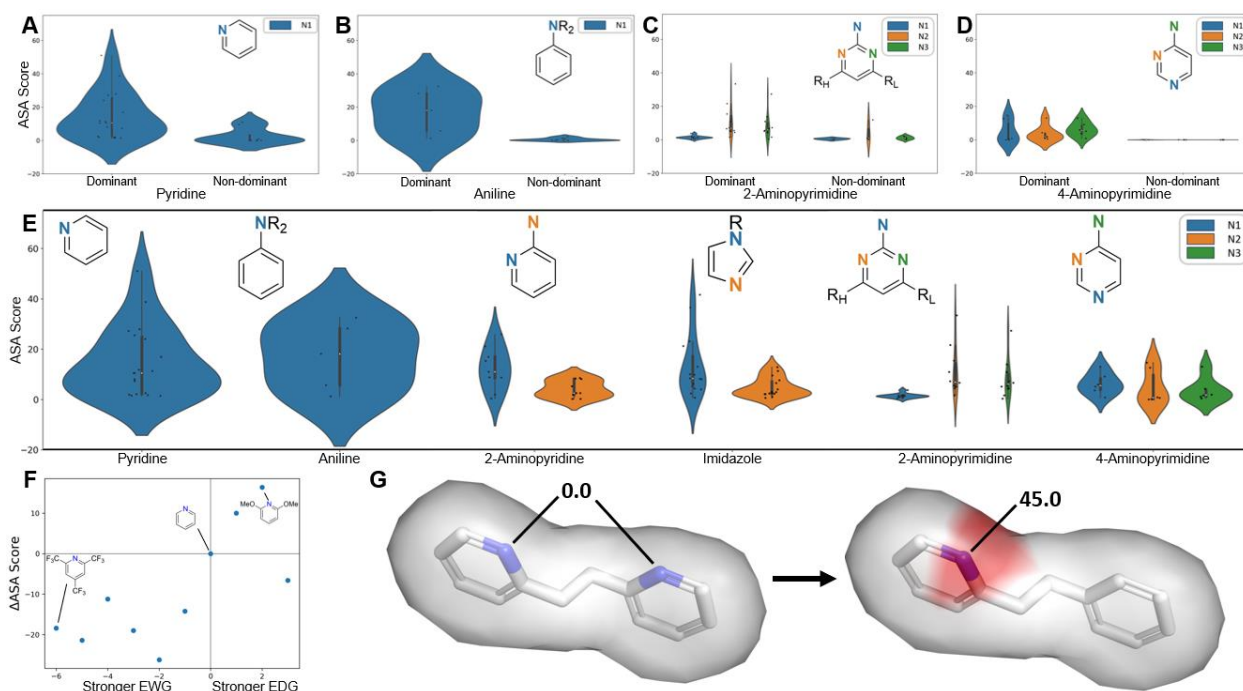
256 **ASA reveals BCL-XpKa implicitly learns substructural information despite substructure-**
257 **free embeddings**

258 To investigate ASA scoring consistency further, we performed ASA on Nitrogen atoms in the
259 most frequently occurring substructures in the basic test set, beyond the controls discussed
260 above. All subgroups examined demonstrated statistically significant loss of ASA signal when a
261 more dominant subgroup (i.e., more relevant to the molecule's basicity) was present (Figure 4A-
262 D), indicating that retaining the most dominant ASA atom also preserves BCL-XpKa's predicted
263 pK_a distribution for the molecule.

264 Filtering out these non-dominant molecules reveals that most substructures have consistent,
265 substructure-specific trends in ASA scores (Figure 4E). This substructure-specificity even
266 persists when separate substructures have Nitrogen atoms with identical local atomic
267 environments. For example, indole and imidazole both have a Nitrogen bound to a Hydrogen and
268 two sp² Carbon atoms, but only imidazole is ionizable at physiologic pH values^{24,25}. While ASA
269 correctly distinguishes that imidazole's other Nitrogen (which has a lone pair of electrons) is the
270 most basic Nitrogen in imidazole, it also scores the N-H Nitrogen significantly higher than the
271 identical motif in indole (Figure 3C), suggesting that this Nitrogen is critical to imidazole's
272 observed basicity.

273 Some of these structures have surprisingly high variance in ASA scores, particularly the pyridine
274 and aniline substructures. From manual inspection, we hypothesized a portion of this variance is
275 attributable to the impact of neighboring electron-donating and -withdrawing groups (EDGs,
276 EWGs), which respectively increase and decrease basicity of neighboring Nitrogen atoms. We
277 tested this hypothesis by scoring manually created sets of pyridine derivatives with various EDG
278 and EWG substituents, which confirmed as suspected that neighboring EDGs tend to increase
279 ASA scores, and neighboring EWGs tend to decrease ASA scores (Figure 4F). Interestingly,
280 symmetric substructures also contributed to this variance, as the symmetric substructure masks
281 the effect of the removed atom during ASA scoring (Figure 4G).

282 Together, these ASA findings suggest that BCL-XpKa has learned impressive substructural
283 insights that are adaptable to molecular context without directly encoding these substructures in
284 the feature set.



285

286 **Figure 4: Atomic sensitivity analysis of substructures** (A-D) Violin plots of ASA scores for
 287 commonly occurring substructures when these substructures are the dominant site of a
 288 molecule's ionization vs when a more dominant substructure was present. (E) Violin plots of
 289 ASA scores for commonly occurring substructures when these substructures were the dominant
 290 site of ionization. Notably, all test-set bases containing 2-Aminopyrimidine and Imidazole featured
 291 them as their dominant ionization site. (F) Change in pyridine Nitrogen's ASA Score by
 292 neighboring EWG or EDG groups. (F) Masking effect of molecular symmetry on ASA score.
 293 ASA = Atomic Sensitivity Analysis; EWG = electron-withdrawing group; EDG = electron-
 294 donating group.

295

296 Atomic sensitivity analysis can inform lead compound optimization in drug development

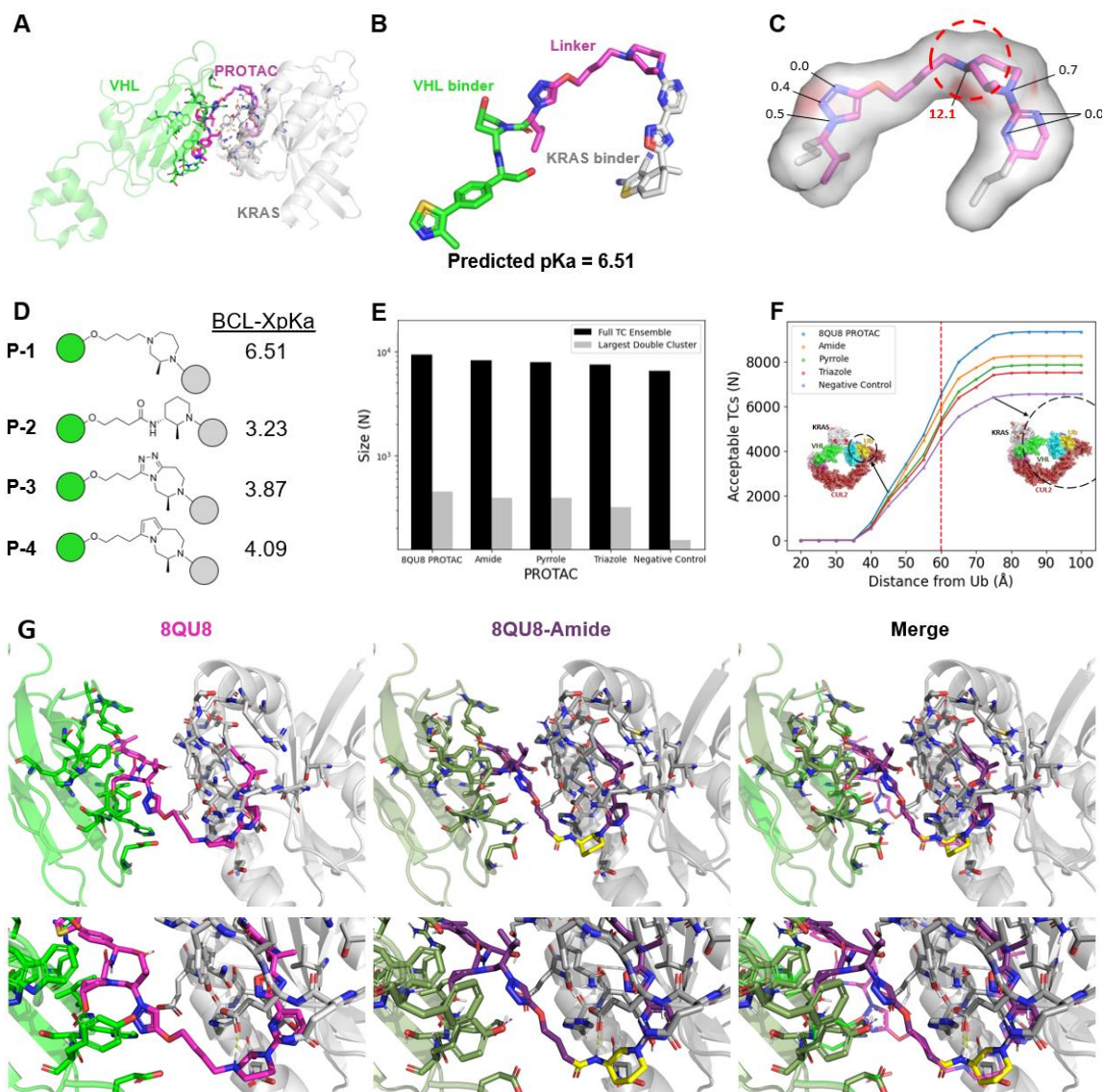
297 Atomic sensitivity analysis also has promising utility in prospective drug design. Significant
 298 interest has developed in the past two decades in targeted protein degradation via small-molecule
 299 Proteolysis Targeting Chimeras (PROTACs). PROTACs consist of two small molecules joined
 300 by a linker and form flexible ternary complexes with the target and an E3 Ligase, which allows
 301 for target ubiquitination and subsequent degradation by the 26S proteasome (Figure 5A, PDB:
 302 8QU8). PROTACs degrade targets catalytically, making them an attractive strategy for
 303 challenging targets that have evaded small-molecule inhibition; however, PROTAC size and
 304 complexity plagues design efforts with poor bioavailability and cell permeability²⁶. These
 305 properties are generally optimized through modification of the PROTAC linker, as the Ligase-
 306 and target-binding domains make specific contacts with their respective proteins. Here, we
 307 demonstrate how atomic sensitivities can guide rational changes to the PROTAC linker to
 308 minimize PROTAC ionizability.

309 KRAS is a challenging target in oncology because it lacks a deep, well-defined pocket for small-
310 molecule inhibition²⁷. Recently, Popow et al. created a KRAS-degrading PROTAC using the
311 VHL E3 Ligase (Figure 5A-B)²⁸. Our model predicts this PROTAC has a pK_a of 6.51 (P-1,
312 Figure 4B) , suggesting protonatability at physiologic pH values. Atomic sensitivity analysis of
313 the linker reveals one of two tertiary amines in P-1 drives this prediction (Figure 5C). The crystal
314 structure of the P-1 ternary complex demonstrates this amine forms a salt bridge with KRAS
315 Q62 when protonated. While salt-bridge interactions promote strong drug binding, PROTACs
316 only need to bind their targets transiently, and a protonatable amine is a liability for permeability.
317 Based on this analysis, we evaluated several P-1 bioisosteric modifications at this amine that
318 reduce predicted pK_a (Figure 5D).

319 We evaluated each modification's performance computationally using a well-benchmarked
320 ternary complex (TC) generation algorithm, which produces a TC ensemble and "double
321 clusters" this ensemble by structural similarity of both the Ligase-Target interactions and
322 PROTAC conformation²⁹. Per the algorithm authors' benchmark, the largest structural double
323 cluster often yields the most crystal-structure-like poses and filters out false positives from the
324 full ensemble. Each modification produced comparable TC ensembles and double clusters to P-1
325 (Figure 5E).

326 Emerging experimental evidence suggests that a PROTAC's degradation efficiency depends on
327 its ability to place the target protein in proximity to ubiquitin (<60Å) in the complete E3 Ligase
328 complex³⁰. As such, we also evaluated the modifications' ability to place KRAS in proximity to
329 ubiquitin by superposing each TC ensemble member into the full E3 Ligase complex. Here, we
330 find that the amide modification in P-2 provides KRAS similar ubiquitin access to P-1 and
331 superior access to all other modifications tested (Figure 5F). Furthermore, docking P-2 into the
332 VHL-KRAS pose identified in the 8QU8 crystal structure demonstrates that P-2 can recapitulate
333 the binding pose of P-1 (heavy atom RMSD <2.5Å), including the key hydrogen bond with
334 KRAS Q62 that P-1 utilizes using the tertiary amine (Figure 5G).

335 Together, these results demonstrate a computationally validated use of atomic sensitivities to
336 guide lead-molecule optimization.



337

338

339 **Figure 5: Atomic sensitivity for drug design** (A) Crystal structure (PDB: 8QU8) of pan-KRAS
 340 degrading PROTAC **P-1** in ternary complex with VHL and KRAS. (B) PROTAC **P-1** colored
 341 according to $5A$, with pK_a calculated by BCL-XpKa. (C) ASA scores for **P-1** Nitrogen atoms.
 342 VHL- and KRAS-binders omitted for space. (D) Proposed bioisosteric **P-1** linker modifications
 343 with pK_a values predicted by BCL-XpKa. (E) Ternary-complex ensemble size and size of largest
 344 Protein-PROTAC conformational double cluster for each linker modification in 5D. (F)
 345 Cumulative number of ternary complexes near Ubiquitin at increasing distances from Ubiquitin
 346 in the closed conformation of the CUL2 E3 Ligase complex. Red line at 60Å denotes an
 347 empirically estimated distance beyond which target ubiquitination is improbable. (G)
 348 Representative images of the 8QU8 crystal structure and the Amide-based linker modification **P-2**
 349 supporting similar PROTAC conformations that preserve the hydrogen bond to KRAS Q62. **P-1**
 350 complex shown in brighter colors; **P-2** complex shown in muted colors; **P-2** linker modification
 351 highlighted in yellow.

352 DISCUSSION

353 Here, we have presented BCL-XpKa, a deep-learning based pK_a predictor that reframes QSPR
354 prediction as a classification problem and avoids explicit substructural embeddings while
355 maintaining competitiveness with contemporary machine learning pK_a predictors. We found that
356 this multitask classification approach directly informs the model's uncertainty in its prediction,
357 and that, beyond its absolute accuracy, BCL-XpKa reliably predicts the effects of common
358 molecular modifications made to a hit/lead compound in a drug development program. We also
359 showed that BCL-XpKa generalizes to foreign substructures better than equivalent models
360 trained on MACCS-based descriptors via leave-substructural-class-out validation.

361 We then used BCL-XpKa as a model system to introduce atomic sensitivity analysis (ASA), a
362 first-in-class ML interpretability method we designed to provide actionable insights into QSPR
363 model output by decomposing a molecule's QSPR prediction into its atomic contributions
364 through direct perturbation of the input chemical structure. When applied to BCL-XpKa, ASA
365 identifies the most ionizable atoms in both acids and bases with remarkable accuracy. ASA also
366 revealed surprisingly consistent results for how BCL-XpKa considers ionizable substructures at
367 the atomic level. These substructural ASA scores were responsive to neighboring electron
368 donating and withdrawing groups, demonstrating that BCL-XpKa learns context-dependent
369 substructural information without explicit substructural embeddings. Finally, we showed that
370 pairing a QSPR model's molecule-level predictions with atomic-level contributions is a powerful
371 tool for guiding lead optimization using a published KRAS-degrading PROTAC. Here, BCL-
372 XpKa and ASA directed linker modifications that reduced PROTAC ionizability while retaining
373 critical PROTAC-KRAS contacts from the original crystal structure.

374 Several limitations exist in our current framework. First, while regression models can predict
375 arbitrarily extreme values given enough quality data, BCL-XpKa must place all extreme values
376 in two catch-all bins, " $pK_a < 0$ " and " $pK_a > 12$ ", given its multitask classifier architecture. While
377 this limits BCL-XpKa's theoretical output range to -0.5 to 12.5, this is not consequential for
378 biologically relevant pH scales and only marginally affects prediction accuracy.

379 Further, ASA is currently limited to atomic-level model explainability. While this provides
380 excellent resolution for atomic properties like pK_a , there are many QSPR tasks where
381 understanding the contribution of entire substructures or pharmacophores would be valuable.
382 Generalizing ASA to higher order molecular substructures will further expand our understanding
383 of QSPR ML model predictions and allow ASA to be tailored to specific tasks.

384 BCL-XpKa and ASA have fundamental applications in computational chemistry generally, as
385 well as early- and late-stage drug development. First, ASA is a generalizable strategy that can
386 increase the explainability of any machine learning model that uses chemical structures as input
387 data. As shown here, ASA scores can help scientists understand what their model has learned
388 from their training set. This information can then guide training-set data augmentation or feature-
389 set modifications.

390 Further, BCL-XpKa paired with ASA is positioned well to support high-quality small-molecule
391 structure preparation for virtual high-throughput screening (vHTS). vHTS involves screening

392 ultra-large libraries (ULLs) of small molecules (currently nearing 10^{11} molecules) for their
393 ability to bind to a protein target. vHTS has notoriously low hit rates, and improper protonation
394 of ULL molecules can contribute to both false-positive and false-negative vHTS screens. BCL-
395 XpKa and ASA's speed and accuracy at predicting pK_a and ionization sites in multiprotic species
396 make this tool a valuable asset for ULL structure preparation and downstream protein-ligand
397 analysis in vHTS.

398 Finally, as demonstrated here, BCL-XpKa paired with ASA can identify ionizable regions in a
399 compound for modification in hit-to-lead or lead optimization. While ionization-site
400 identification is relatively straightforward, this model-ASA strategy generalizes to any
401 QSPR/QSAR model. Therefore, applying ASA to predictors of ADMET/DMPK may facilitate
402 understanding of important but less readily interpretable liabilities in a hit or lead compound.

403

404 **METHODS**

405 *Training Datasets*

406 ChEMBL27 is an open-source database contains over 2 million molecules with various
407 physicochemical descriptors³¹. ACDlabs was used to calculate acidic pK_a and basic pK_a values
408 (chembl_acid_pka and chembl_base_pka, respectively) for ChEMBL molecules, and molecules
409 that were included in our test sets were excluded¹⁸. We also generated negative data (molecules
410 with no ionization site) in the BCL and set chembl_acid_pka = 50, chembl_base_pka = 0. In
411 sum, acidic pK_a models were trained on 988,643 molecules, and basic pK_a models on 812,918
412 molecules.

413

414 *Molecule Preparation*

415 Molecular 3D structures were standardized using Corina for training and testing BCL-XpKa³².
416 For external models, structure preparation followed the authors' direction, and Corina was used
417 if no structure preparation method was mentioned. This standardization was used exclusively for
418 downstream usability, as BCL-XpKa solely uses 2D descriptors. All PROTAC modifications
419 introduced in Figure 5 were minimized in the Molecular Operating Environment (MOE) prior to
420 ternary complex ensemble generation.

421

422 *Molecular Features*

423 The Mol2D molecular descriptor set was used to encode molecules as described elsewhere¹⁷.
424 Briefly, for each atom in a molecule, Mol2D encodes information about that atom, the bonds
425 made to that atom, and the atoms one bond length away from that atom.

426 To train a multitask classifier with N output labels, ChEMBL pK_a values were encoded in an
427 $N \times 1$ result label $R \in \mathbb{Z}_2^N$, where the first entry corresponds to $pK_a < 0$, the last entry to $pK_a \geq$

428 12, and the i^{th} entry to $(12/N) * (i - 1) \leq \text{pK}_a < (12/N) * i$. For regression models, the
429 pK_a in the ChEMBL set was used directly as the result label.

430

431 *Model Training and Validation*

432 Artificial neural networks were built in C++ for the Biology and Chemistry Library (BCL), an
433 open-source cheminformatics platform created and maintained by our lab. Each model was
434 trained for 250 iterations without early stopping, which our lab previously found to be
435 unnecessary when dropout is used³³. An upper-bound for model performance was calculated
436 through random-split cross validation. A lower-bound for performance was calculated through
437 leave-class-out cross validation (LCO-CV), in which the training set was divided into 30 subset
438 $\{C_i\}_{i=1}^{30}$ based on ionizable groups defined in literature³⁴. 30 models were then trained in an
439 iterative all-but-one scheme. Model internal performance was evaluated using the logarithmic
440 receiver operating characteristic curve (logAUC) and AUC using the BCL
441 model:ComputeStatistics application.

442

443 *Model Output and Evaluation*

444 MTCs with N output labels calculate N local positive predictive values (localPPV), where the i^{th}
445 localPPV denotes the probability that the pK_a lies in the i^{th} pK_a interval (see *Molecular Features*
446 above). For model evaluation, we report mean absolute error (MAE) for reader familiarity. In the
447 supplement, we also provide a Brier score for each model, which is a proper scoring rule¹ for
448 more rigorous evaluation of classification model output.

449 For each molecule $m_i \in Y$ in a test set Y with M molecules, the pK_a of m_i was encoded with a
450 binary result label $R_i \in \mathbb{Z}_2^N$ as described above. MTC scored each molecule, providing a discrete
451 probability distribution P_i describing the molecule's likely pK_a interval membership. From these
452 distributions, MAE was calculated as:

$$453 \quad \text{MAE}(P, Y) = \frac{1}{M} \sum_{i=1}^M |y_i - E[P_i]|$$

454 where $E[P_i]$ is the expected value of P_i . Similarly, for several models a Brier score was
455 calculated as:

456

$$457 \quad \text{BS}(P, Y) = \frac{1}{M} \sum_{i=1}^M \sum_{j=1}^N (r_{ij} - p_{ij})^2$$

¹ "Proper scoring rule" is a term in statistics for a loss function that is minimized if the probability distribution output from the model is identical to the ground-truth probability distribution. When the bidirectional holds, the scoring rule is further labeled a *strictly proper* scoring rule.

458 Where $r_{ij} \in R_i$ is the j^{th} result label for the i^{th} molecule in the test set (i.e., whether its pK_a value
459 lies in the j^{th} pK_a interval), and p_{ij} is the localPPV that the i^{th} molecule's pK_a value lies in the j^{th}
460 pK_a interval.

461 Throughout, MAE is used throughout to compare MTC to Regression and MTC to MTC models.
462 Percent accuracy of categorization is not included, as it is an improper and discontinuous scoring
463 metric.

464

465 *Atomic sensitivity analysis*

466 Atom replacement schemes were coded in C++ within the BCL. A parent molecule m is scored
467 by BCL-XpKa to produce P , a discrete probability distribution of potential pK_a values.
468 Heteroatom a in the parent molecule is replaced with an appropriately hybridized carbon atom,
469 and the perturbed molecule is rescored to produce P'_a . The dissimilarity between these
470 distributions was calculated by their Kullback-Leibler (KL) divergence:

$$471 D_{\text{KL}}(P'_a||P) = \sum_{j=1}^N P'_{a,j} \ln \left(\frac{P'_{a,j}}{P_j} \right),$$

472 where P_j and $P'_{a,j}$ are localPPVs as described in *Model Output and Evaluation*. Briefly, the KL
473 divergence of these two probability distributions is best interpreted as the relative entropy
474 between these distributions, where $D_{\text{KL}}(P'||P) = 0$ denotes the distributions are identical (there
475 would be no “surprises” if a given sample came from P vs P'), and higher values denote more
476 dissimilarity. Finally, KL divergences were empirically denoised to generate ASA scores:

$$477 \text{ASA}(m, a) = e^{\lfloor 5 * D_{\text{KL}}(P'_a||P) \rfloor} - 1$$

478

479 *PROTAC ternary complex ensemble generation*

480 Ternary complexes (TCs) were constructed according to Drummond et al (2020). Briefly,
481 protein-protein interactions (PPIs) with the PROTAC binding pockets near each other, as well as
482 a set of up to 10000 PROTAC conformations, were produced in the Molecular Operating
483 Environment (MOE). PROTAC conformations were then docked into the PPIs and filtered
484 according to the authors' criteria. TCs were then clustered on both protein- and PROTAC-
485 conformations to produce “double clusters.” Protein-conformational clustering was done at CA-
486 RMSD < 10Å. PROTAC clustering was done at heavy-atom RMSD < 2.5Å.

487

488 *Hardware*

489 All models were trained with 18 Intel Xenon W-2295 CPU cores. PROTAC TC formation was
490 performed using an Nvidia RTX A5000 GPU.

491

492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537

REFERENCES

- 1 Gorgulla, C. Recent Developments in Ultralarge and Structure-Based Virtual Screening Approaches. *Annu Rev Biomed Data Sci* **6**, 229-258, doi:10.1146/annurev-biodatasci-020222-025013 (2023).
- 2 Dara, S., Dhamercherla, S., Jadav, S. S., Babu, C. M. & Ahsan, M. J. Machine Learning in Drug Discovery: A Review. *Artif Intell Rev* **55**, 1947-1999, doi:10.1007/s10462-021-10058-4 (2022).
- 3 Maltarollo, V. G., Kronenberger, T., Wrenger, C. & Honorio, K. M. Current Trends in Quantitative Structure–Activity Relationship Validation and Applications On Drug Discovery. *Future Science OA* **3**, FSO214, doi:10.4155/fsoa-2017-0052 (2017).
- 4 Wu, J., Kang, Y., Pan, P. & Hou, T. Machine learning methods for pKa prediction of small molecules: Advances and challenges. *Drug Discovery Today* **27**, 103372, doi:<https://doi.org/10.1016/j.drudis.2022.103372> (2022).
- 5 Navo, C. D. & Jiménez-Osés, G. Computer Prediction of pKa Values in Small Molecules and Proteins. *ACS Medicinal Chemistry Letters* **12**, 1624-1628, doi:10.1021/acsmchemlett.1c00435 (2021).
- 6 Johnston, R. C. *et al.* Epik: pKa and Protonation State Prediction through Machine Learning. *Journal of Chemical Theory and Computation* **19**, 2380-2388, doi:10.1021/acs.jctc.3c00044 (2023).
- 7 Pan, X., Wang, H., Li, C., Zhang, J. Z. H. & Ji, C. MolGpka: A Web Server for Small Molecule pK(a) Prediction Using a Graph-Convolutional Neural Network. *J Chem Inf Model* **61**, 3159-3165, doi:10.1021/acs.jcim.1c00075 (2021).
- 8 Xiong, J. *et al.* Multi-instance learning of graph neural networks for aqueous pKa prediction. *Bioinformatics* **38**, 792-798, doi:10.1093/bioinformatics/btab714 (2021).
- 9 Gierlich, C. & Palkovits, S. Featurizing chemistry for machine learning — methods and a coded example. *Current Opinion in Chemical Engineering* **37**, 100840, doi:<https://doi.org/10.1016/j.coche.2022.100840> (2022).
- 10 BALLS, G. R., PALMER-BROWN, D. & SANDERS, G. E. Investigating microclimatic influences on ozone injury in clover (*Trifolium subterraneum*) using artificial neural networks. *New Phytologist* **132**, 271-280, doi:<https://doi.org/10.1111/j.1469-8137.1996.tb01846.x> (1996).
- 11 Maier, H. R. & Dandy, G. C. The Use of Artificial Neural Networks for the Prediction of Water Quality Parameters. *Water Resources Research* **32**, 1013-1022, doi:<https://doi.org/10.1029/96WR03529> (1996).
- 12 Štrumbelj, E., Kononenko, I. & Robnik Šikonja, M. Explaining instance classifications with interactions of subsets of feature values. *Data & Knowledge Engineering* **68**, 886-904, doi:<https://doi.org/10.1016/j.datak.2009.01.004> (2009).
- 13 Fong, R. C. & Vedaldi, A. Interpretable Explanations of Black Boxes by Meaningful Perturbation. *2017 IEEE International Conference on Computer Vision (ICCV)*, 3449-3457 (2017).
- 14 McCloskey, K., Taly, A., Monti, F., Brenner, M. P. & Colwell, L. J. Using attribution to decode binding mechanism in neural network models for chemistry. *Proc Natl Acad Sci U S A* **116**, 11624-11629, doi:10.1073/pnas.1820657116 (2019).
- 15 Olden, J. D. & Jackson, D. A. Illuminating the “black box”: a randomization approach for understanding variable contributions in artificial neural networks. *Ecological Modelling* **154**, 135-150, doi:[https://doi.org/10.1016/S0304-3800\(02\)00064-9](https://doi.org/10.1016/S0304-3800(02)00064-9) (2002).
- 16 Karpov, P., Godin, G. & Tetko, I. V. Transformer-CNN: Swiss knife for QSAR modeling and interpretation. *Journal of Cheminformatics* **12**, 17, doi:10.1186/s13321-020-00423-w (2020).

538 17 Vu, O., Mendenhall, J., Altarawy, D. & Meiler, J. BCL::Mol2D—a robust atom environment
539 descriptor for QSAR modeling and lead optimization. *J Comput Aided Mol Des* **33**, 477–486,
540 doi:10.1007/s10822-019-00199-8 (2019).

541 18 Percepta-Batch (Advanced Chemistry Development, Inc. (ACD/Labs), Toronto, ON, Canada,
542 2023).

543 19 Marvin (ChemAxon, Ltd, 2023).

544 20 Molecular Operating Environment (MOE) (Chemical Computing Group, Inc., Montreal, Quebec,
545 Canada, 2023).

546 21 Mansouri, K., Grulke, C. M., Judson, R. S. & Williams, A. J. OPERA models for predicting
547 physicochemical properties and environmental fate endpoints. *J Cheminform* **10**, 10,
548 doi:10.1186/s13321-018-0263-1 (2018).

549 22 de Souza Neto, L. R. *et al.* In silico Strategies to Support Fragment-to-Lead Optimization in Drug
550 Discovery. *Frontiers in Chemistry* **8**, doi:10.3389/fchem.2020.00093 (2020).

551 23 Jiménez-Luna, J., Grisoni, F. & Schneider, G. Drug discovery with explainable artificial
552 intelligence. *Nature Machine Intelligence* **2**, 573–584, doi:10.1038/s42256-020-00236-4 (2020).

553 24 Foulon, C. *et al.* Determination of ionization constants of N-imidazole derivatives, aromatase
554 inhibitors, using capillary electrophoresis and influence of substituents on pKa shifts. *Journal of*
555 *Chromatography A* **1035**, 131–136, doi:<https://doi.org/10.1016/j.chroma.2004.02.053> (2004).

556 25 Hinman, R. L. & Lang, J. The Protonation of Indoles. Basicity Studies. The Dependence of Acidity
557 Functions on Indicator Structure. *Journal of the American Chemical Society* **86**, 3796–3806,
558 doi:10.1021/ja01072a040 (1964).

559 26 Hornberger, K. R. & Araujo, E. M. V. Physicochemical Property Determinants of Oral Absorption
560 for PROTAC Protein Degraders. *Journal of Medicinal Chemistry* **66**, 8281–8287,
561 doi:10.1021/acs.jmedchem.3c00740 (2023).

562 27 Wu, X. *et al.* Small molecular inhibitors for KRAS-mutant cancers. *Frontiers in Immunology* **14**,
563 doi:10.3389/fimmu.2023.1223433 (2023).

564 28 Popow, J. *et al.* Targeting cancer with small molecule pan-KRAS degraders. *bioRxiv*,
565 2023.2010.2024.563163, doi:10.1101/2023.10.24.563163 (2023).

566 29 Drummond, M. L., Henry, A., Li, H. & Williams, C. I. Improved Accuracy for Modeling PROTAC-
567 Mediated Ternary Complex Formation and Targeted Protein Degradation via New In Silico
568 Methodologies. *J Chem Inf Model* **60**, 5234–5254, doi:10.1021/acs.jcim.0c00897 (2020).

569 30 Dixon, T. *et al.* Predicting the structural basis of targeted protein degradation by integrating
570 molecular dynamics simulations with structural mass spectrometry. *Nature Communications* **13**,
571 5884, doi:10.1038/s41467-022-33575-4 (2022).

572 31 Zdrazil, B. *et al.* The ChEMBL Database in 2023: a drug discovery platform spanning multiple
573 bioactivity data types and time periods. *Nucleic Acids Research* **52**, D1180–D1192,
574 doi:10.1093/nar/gkad1004 (2023).

575 32 Sadowski, J. & Gasteiger, J. From atoms and bonds to three-dimensional atomic coordinates:
576 automatic model builders. *Chemical Reviews* **93**, 2567–2581, doi:10.1021/cr00023a012 (1993).

577 33 Mendenhall, J. & Meiler, J. Improving quantitative structure-activity relationship models using
578 Artificial Neural Networks trained with dropout. *J Comput Aided Mol Des* **30**, 177–189,
579 doi:10.1007/s10822-016-9895-2 (2016).

580 34 Ropp, P. J., Kaminsky, J. C., Yablonski, S. & Durrant, J. D. Dimorphite-DL: an open-source program
581 for enumerating the ionization states of drug-like small molecules. *Journal of Cheminformatics*
582 **11**, 14, doi:10.1186/s13321-019-0336-9 (2019).

583