

# Learning From The Ligand: Using Ligand-Based Features To Improve Binding Affinity Prediction

Fergus Boyles, Charlotte M. Deane, and Garrett M. Morris\*

*Department of Statistics, University of Oxford, 24-29 St Giles', Oxford, OX1 3LB, U.K.*

E-mail: [morris@stats.ox.ac.uk](mailto:morris@stats.ox.ac.uk)

Phone: +44 1865 281770. Fax: +44 1865 282862

## Abstract

Machine learning scoring functions for protein-ligand binding affinity prediction have been found to consistently outperform classical scoring functions. Structure-based scoring functions for universal affinity prediction typically use features describing interactions derived from the protein-ligand complex, with limited information about the chemical or topological properties of the ligand itself. We demonstrate that the performance of machine learning scoring functions are consistently improved by the inclusion of diverse ligand-based features. For example, a Random Forest combining the features of RF-Score v3 with RDKit molecular descriptors achieved Pearson correlation coefficients of up to 0.831, 0.785, and 0.821 on the PDBbind 2007, 2013, and 2016 core sets respectively, compared to 0.790, 0.737, and 0.797 when using the features of RF-Score v3 alone. Excluding proteins and/or ligands that are similar to those in the test sets from the training set has a significant effect on scoring function performance, but does not remove the predictive power of ligand-based features. Furthermore a Random Forest using only ligand-based features is predictive at a level similar to classical scoring functions and it appears to be predicting the mean binding affinity of a

ligand for its protein targets. Data and code to reproduce all results freely available at <http://opig.stats.ox.ac.uk/resources>.

## Introduction

Structure-based virtual screening (SBVS) uses the 3D structure of a target protein to screen large compound libraries for small molecules likely to bind. “Explicit” SBVS uses protein-ligand docking to predict the binding mode of each compound within the active site, and a scoring function to predict the strength of binding<sup>1-5</sup>. While it is possible to compute the binding affinity of a compound using more rigorous methods such as free energy perturbation<sup>6</sup>, their computational cost makes them impractical for screening libraries of millions of compounds<sup>7</sup>. To save time, in SBVS, more approximate scoring functions designed to estimate the binding affinity are used. Scoring functions are typically employed for three tasks in SBVS: correctly identifying the binding mode of a ligand (pose prediction or “docking”); classifying molecules as either active or inactive (“virtual screening”); and ranking ligands in order of their binding affinity for a given protein target (“scoring”). Popular protein-ligand docking packages, such as GOLD<sup>8</sup>, Glide<sup>9,10</sup>, ICM<sup>11</sup>, FlexX<sup>12</sup>, Surflex<sup>13</sup>, and the AutoDock family<sup>14-16</sup>, rely on a single scoring function to perform all three tasks simultaneously. These scoring functions make use of molecular force fields, statistical potentials, or linear combinations of empirical terms to assign a score to a receptor-ligand complex, and are often collectively referred to as ‘classical’ scoring functions. While docking methods have been used successfully to predict binding modes, and for virtual screening, correctly ranking ligands by their binding affinity for a protein remains extremely challenging<sup>17,18</sup>.

Recently, however, the application of machine learning techniques has led to the development of new scoring functions that outperform classical scoring functions in terms of ranking compounds by binding affinity<sup>19-23</sup>. Like classical scoring functions, these methods use the 3D structure of the protein and an automatically generated binding

mode of the ligand to compute structure-based interaction features between the protein and ligand. The binding mode-derived features are then used as inputs to a machine learning algorithm to predict the binding affinity. The features used by such machine learning scoring functions typically focus on capturing interactions between the ligand and the protein, but make limited use of the bulk physical, chemical, and topological properties of the protein and/or the ligand alone. Ligand-based features are widely used to select potential binders in ligand-based virtual screening, and have proven to be very effective in understanding polypharmacological relationships between proteins<sup>24</sup>. Although scoring functions such as SFCscore<sup>21,25</sup>, NNScore 2.0<sup>20</sup>, and those used in AutoDock 4<sup>14</sup> and AutoDock Vina<sup>15</sup> include some features of the ligand, the use of detailed information about the ligand to predict cognate protein-ligand binding affinity remains limited. Perhaps the closest is the field of proteochemometric modeling,<sup>26</sup> in which classification models are constructed using a combination of features describing the protein and the small molecule. This approach has proven to be capable of predicting ligand selectivity<sup>27</sup> and capturing polypharmacology.<sup>28</sup>

Motivated by the utility of ligand-based features in virtual screening and proteochemometric modelling, we investigated whether a more detailed representation of the ligand can improve the ability of a scoring function to predict its binding affinity. Using the cheminformatics toolkit RDKit (<https://www.rdkit.org/>, accessed 17/05/2019) we computed a diverse set of 1D and 2D ligand molecular descriptors and combined these with the structure-based features used by the machine learning scoring functions RF-Score<sup>19</sup>, RF-Score v3<sup>29</sup>, NNScore 2.0<sup>20</sup>, as well as the empirical scoring function of AutoDock Vina<sup>15</sup>. We show that a Random Forest (RF)<sup>30</sup> regression model using both structure-based and ligand-based features consistently outperforms a model using only structure-based features when benchmarked on three versions of the PDBbind<sup>31</sup> core set, corresponding to the scoring power test of the three Comparative Assessment of Scoring Functions<sup>17,18,32,33</sup>.

We find that removal of test-set similar proteins and ligands from the training sets degrades performance but does not abrogate the predictive power of ligand-based

features. Furthermore, we show that a model using only ligand-based features appears to be predictive of the mean affinity of a ligand for its binding partners when trained and tested on PDBbind data. We computed the relative importance of the features used by each model and show that when structure-based and ligand-based features are combined, both structure-based and ligand-based features are among the top-ranked features. These results suggest that quickly-computed ligand-based features should be used to improve the ability of a machine learning scoring function to predict protein-ligand binding affinity.

## Methods

### Training and Test Sets

The PDBbind database<sup>31</sup> is a curated set of bound macromolecule structures drawn from the Protein Data Bank (PDB)<sup>34</sup>, each with an experimentally-measured binding affinity for its binding partner. Each release of PDBbind includes a “general set”, which contains all the protein-ligand structures in the database; and a “refined set”, a subset of protein-ligand complexes satisfying strict criteria concerning structure quality, affinity data reliability, and the nature of the complex. The 2018 release of PDBbind contains 16,151 protein-ligand complexes in the general set, with 4,463 complexes in the refined set. We used the refined set as our primary source of training data; however, it has been reported that including the lower-quality data comprising the remainder of the general set can still improve the performance of machine learning scoring functions<sup>35</sup>, so we repeated our analysis using the general set as our source of training data.

To validate our models, we used a subset of the PDBbind refined set referred to as the “core set”. This is obtained by clustering the proteins in the refined set at 90% sequence identity and selecting three or five (depending on the version of PDBbind) representatives of each cluster for which the corresponding ligands have a broad range of binding affinity values, resulting in a diverse, non-redundant set of protein-ligand

complexes. We repeated our tests using the core sets from the 2007, 2013, and 2016 releases of PDBbind. These versions of the core set were used as the “scoring power” benchmark in the Comparative Assessment of Scoring Functions (CASF) exercises: CASF2009<sup>32</sup>, CASF2013<sup>17,18</sup>, and CASF2016<sup>33</sup> respectively, allowing our results to be compared directly to previously published scoring function benchmarks. We used the test sets corresponding to all three CASF benchmarks since the contents of each test set are substantially different to the others and so each set offers a different challenge for a scoring function. Excluding proteins and ligands that could not be parsed by OpenBabel or RDKit resulted in 2007, 2013, and 2016 core sets containing 196, 180, and 276 structures respectively. The full list of structures omitted from the core sets because of parsing failures is included in the Supporting Information. These test sets are relatively small, making it difficult to identify statistically-significant differences in results generated by different models. To partially overcome this shortcoming without deviating from widely-used benchmarks, we combined the structures from each core set into a fourth test set with duplicate structures removed. This combined test set numbered 525 structures, almost twice the size of the PDBbind 2016 core set.

PDBbind provides for each complex an experimentally-determined value of the inhibition constant  $K_i$ , the dissociation constant  $K_d$ , or the half-maximal inhibitory concentration  $IC_{50}$ , in decreasing order of preference (*e.g.* if both  $K_i$  and  $K_d$  values are available, PDBbind reports the measurement of  $K_i$ ). The refined set includes only measurements of  $K_i$  and  $K_d$ , while the general set also includes data for which only  $IC_{50}$  measurements were available. For our purposes, these values are used interchangeably and are collectively denoted by the binding constant,  $K$ . We used the negative base-10 logarithm of  $K$ , commonly denoted as  $pK$ :

$$pK = -\log_{10} K$$

We evaluated each scoring function by computing the Pearson correlation coefficient,  $\rho_p$ , between its predictions  $\{\hat{y}\}$  and the experimental values  $\{y\}$  of  $pK$  for the

complexes in the test set. Confidence intervals were estimated at the 95% significance level using 10,000 bootstrapped samples of the predictions. The Mann-Whitney U test was used to compare the distribution of bootstrapped  $\rho_p$  values to assess the significance of the differences in correlation coefficient between two scoring functions. We also performed a permutation test with 10,000 samples to assess the possibility that correlations arose by random chance.

## Ligand-Based Features

To represent the ligand in our models we used a diverse set of molecular descriptors computed using the cheminformatics toolkit RDKit. Using the *Descriptors* module of the Python RDKit package version 2018.03, we computed a set of 200 molecular descriptors for each ligand. These descriptors are conformation-independent and may be categorized as either (computed) experimental properties (*e.g.* molar refractivity, logP) or theoretical descriptors derived from a symbolic representation of the molecule. The theoretical descriptors may be further categorized according to the dimensionality of the representation of the molecule from which they are derived. The conformer-independent descriptors we consider are either 1-D compositional properties (*e.g.* heavy atom counts, bonds counts, and molecular weight) or 2-D topological properties (*e.g.* fragment counts, topological polar surface area, and connectivity index). Any features with zero variance across the data set, or that were null-valued (*i.e.* infinite or not computable) within the data set were excluded. We removed the Ipc index (an information theory-derived descriptor) as it produced extreme numerical values for larger molecules (too large to be represented as 32-bit floats). In total, 185 RDKit descriptors were retained. We refer to this set of ligand-based features as “RDKit descriptors” throughout this work.

## Structure-Based Features

To investigate the effects of augmentation with ligand molecular descriptors, we considered the features of several publicly-available machine learning scoring functions, namely RF-Score<sup>19</sup>, RF-Score v3<sup>29</sup>, and NNScore 2.0<sup>20</sup>. Both RF-Score v3 and NNScore 2.0 include the six terms used by the AutoDock Vina scoring function<sup>15</sup>. We therefore considered the AutoDock Vina terms separately to examine the effect of combining ligand molecular descriptors with just the terms used by a classical empirical scoring function. We computed the features of each of these scoring functions using the implementations provided by the Open Drug Discovery Toolkit (ODDT) version 0.6<sup>36</sup>.

## Varying Training Set Size and Composition

Previous works by numerous authors have demonstrated that both the size of the training set,<sup>35</sup> and similarity between training and test set structures,<sup>37-39</sup> can influence scoring function performance using the PDBbind core set. We investigated the effect of three factors in training set composition on the performance of our models: training set size; similarity of ligands between training and test examples; and similarity of proteins between training examples.

To examine the effect of training set size, we simulated the effect of adding more structural and affinity training data over time by restricting the training set to annual releases of the PDBbind database from 2013 through to 2018. Each release contains more data than the previous releases, so this results in six training sets of increasing size. By training separately on the general and refined sets of each year, we explore two different scenarios: a larger data set of varying quality, and a smaller data set with strict quality controls, giving a total of twelve distinct training scenarios.

To investigate the effect of including similar ligands in both the training and test sets, we used RDKit to compute the Tanimoto similarity between the Morgan fingerprints (radius 2 and 2048 bits) of each pair of ligands. We then constructed a new training set by removing from the available training data any structure whose

ligand had a Tanimoto similarity of greater than or equal to 0.9 to any ligand in the test set.

To study the effect of including similar proteins in the training and test sets, for each version of PDBbind we constructed a series of training sets by removing from the original training set any structures with a protein sequence identity to any protein in the test set above a threshold of sequence identity. Clustering of the entire PDB using BLASTclust at sequence identity values from 30% to 100% computed by BLASTclust were downloaded from the PDB website (<http://www.rcsb.org/pdb/statistics/clusterStatistics.do>, accessed 13/05/2019). Finally, we construct an additional series of training sets by removing both structures with protein sequence identity above the cutoff value, and structures with ligand Tanimoto similarity greater than or equal to 0.9 to any structure in the test set.

When excluding test-set similar complexes from the training data, we treated each test set separately. For example, when testing on the PDBbind 2016 core set, only proteins similar to those found in the 2016 core set were excluded from the training set.

## Scoring Function Construction

For each of the four scoring functions considered (AutoDock Vina, RF-Score, RF-Score v3, and NNScore 2.0) we constructed two sets of features. The first used only the original features of the scoring functions, while the second used the original features of the scoring function, plus the 183 RDKit descriptors of the ligand. Since AutoDock Vina (and hence RF-Score v3 and NNScore 2.0) already use the number of rotatable bonds of the ligand, we dropped this from the set of RDKit descriptors added to avoid including the same feature twice. Finally, to examine whether there is any signal in the RDKit descriptors independent of the structure-based features, we constructed models using only the RDKit descriptors as a separate feature set, resulting in a total of nine different sets of features.

For each set of features and each training set, we built a scoring function by using

Random Forest<sup>30</sup> (RF) regression to fit an estimator for the  $pK$  of a protein-ligand complex. We used the implementation of RF in the Python machine learning library scikit-learn<sup>40</sup>. Although RF is generally robust with respect to hyperparameter choice, we tested the effect of varying the number of trees in the forest ( $n\_estimators$ ) and the maximum number of features considered at each split ( $max\_features$ ). We chose to set  $n\_estimators=500$  and  $max\_features=0.33$  as these values yielded optimal out-of-bag performance (Supporting Information Fig. 2-3). We tested several other ML algorithms but found that RF consistently achieved the best cross-validation scores (Supporting Information Fig. 1; 6; 7).

## Investigating Affinity Predictions for Ligands Found in Multiple Structures

We investigated the affinity predictions for ligands that are found in multiple structures in the PDBbind database. We clustered the structures of the PDBbind 2018 general set by the three-character chemical ID of the ligand as specified in the PDB, and selected all ligands found in at least three structures. Holding out each ligand in turn as a test case, we trained a RF regression model using only the RDKit descriptors on all structures not containing that ligand, and used the resulting model to predict the affinity of that ligand.

We repeated the above process for each ligand by removing from the training set all ligands with a Tanimoto similarity to any test set ligand above a defined threshold. We then trained a new RF regression model and predicted the affinity of the test ligand. The Tanimoto similarity threshold in this process was reduced in steps of 0.1 from 0.9 to 0.1 inclusive.

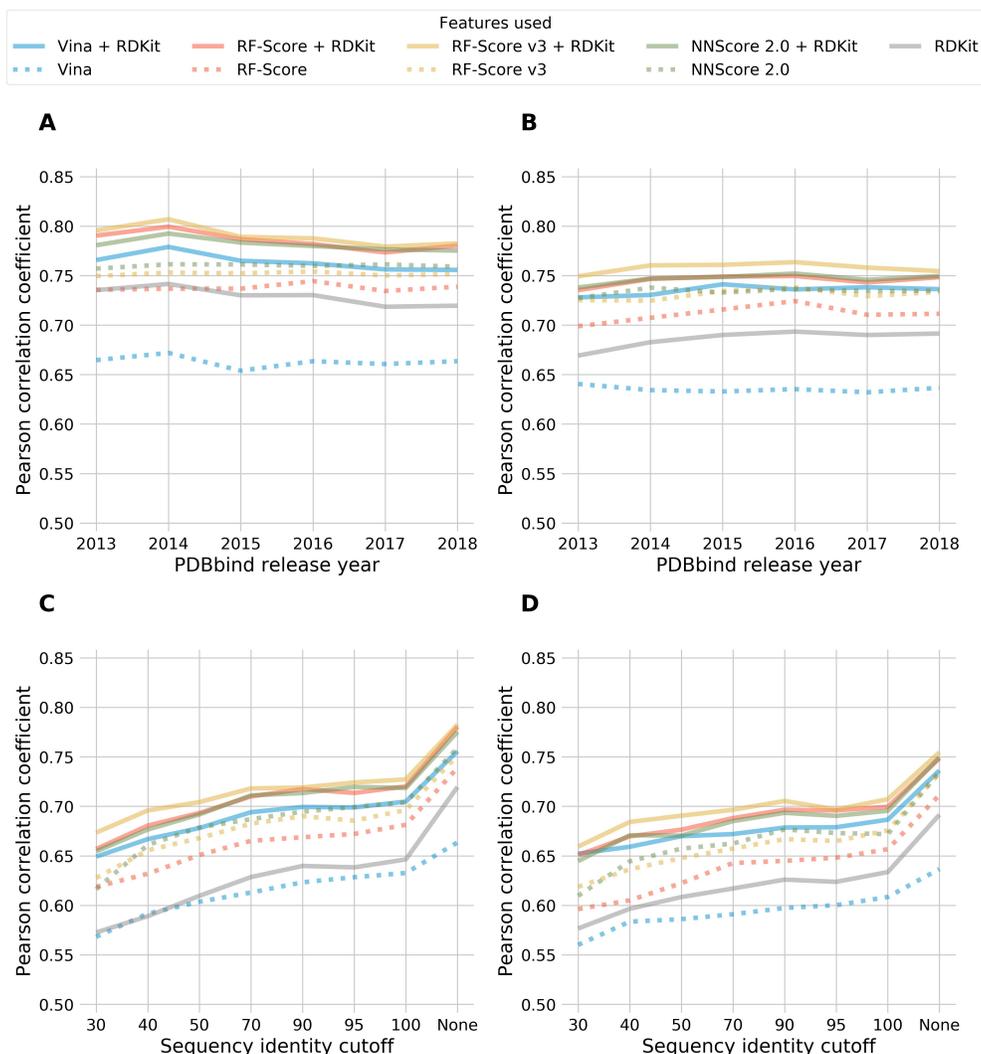


Figure 1: Pearson correlation coefficient for scoring function predictions against experimental pK values for the combined core set. (A, B): Varying training set size chronologically; the version of the PDBbind refined set is indicated on the  $x$ -axis. (C, D): Varying training set composition by removing structures with high protein sequence identity to any test set proteins from the PDBbind 2018 refined set. The sequence identity threshold above which proteins were excluded is indicated on the  $x$ -axis. (A, C): No data excluded on the basis of ligand similarity. (B, D): Structures whose ligand has a Tanimoto similarity of 0.9 or greater to any ligand in the test set were excluded. Solid lines denote ML scoring functions; dotted lines denote ML scoring functions augmented using the RDKit descriptors. The dotted grey line indicates a model using only the RDKit descriptors. The size of the training set has little effect on performance (A, B), while the exclusion of test-set similar proteins from the training set reduces the performance of all SFs (C, D). The exclusion of test-set similar ligands from the training data reduces the performance of all SFs (B, D). In all training scenarios, the SF augmented with ligand-based features outperforms the corresponding structure-based SF (solid line always above the dotted line of the same colour).

# Results and Discussion

## Ligand-Based Features Improve ML Scoring Functions

Fig. 1 shows the Pearson correlation coefficient between experimental and predicted  $pK$  on the combined core set for our nine different scoring functions using four different approaches to training set construction. In all four cases, regardless of training set construction the addition of ligand-based features to a structure-based RF scoring function improves performance (in Fig. 1, for each colour, the solid line showing ligand-based plus structure-based features is consistently above the corresponding dotted line showing structure-based features alone). The trend toward a ligand-feature augmented scoring function outperforming the corresponding scoring function, using only structure-based features alone, is exemplified by the combination of AutoDock Vina terms and RDKit descriptors. This combination results in predictive performance comparable to that of the RF-Score, RF-Score v3, and NNScore 2.0 features, suggesting that a more complex and detailed set of features describing protein-ligand interactions is not necessarily more predictive than a comparatively simple set of force-field-like terms (AutoDock Vina terms) and molecular descriptors of the ligand (RDKit descriptors).

The performance of the scoring functions using each feature set when trained on different versions of the PDBbind refined set with no data removed is shown in Fig. 1A, while Fig. 1B shows corresponding results when structures whose ligands have a high Tanimoto similarity to any ligand in the test set are removed from the training set. In both of these cases there is no consistent improvement in performance when larger, more recent versions of the refined set are used. This is contrary to the results of Li *et al.*<sup>35</sup> and may suggest that, at least when using RF to train an estimator, there is an element of ‘learning saturation’, beyond which there is limited benefit to using additional training data. However, for a given release of PDBbind, a scoring function trained on the general set outperforms the same scoring function trained on the refined set (Supporting Information Fig. 12-13), consistent with the findings of Li *et al.*. This

difference in performance vanishes for all scoring functions when structures with test-set similar ligands are excluded from the training set, and is greatly reduced when structures with 90% protein sequence identity to those in the test set are excluded. This suggests that the increase in performance when training on the general set can be attributed to increased representation of the core set proteins and ligands in the training data.

Exclusion from the training set of test-set similar proteins results in significantly reduced scoring function performance (Fig. 1C). There is a significant drop in the performance of all scoring functions even when a sequence identity cutoff of 100% is imposed, *i.e.* when only proteins with identical sequence to those in the test set are excluded from the training set. Reducing the sequence identity cutoff from 90% to 50% has a smaller impact on performance than the initial imposition of a 100% cutoff. Further reducing the cutoff from 50% to 30% has a more apparent effect. Fig. 1D shows similar behaviour when structures with test-set similar ligands are also excluded from the training set. As in the case where similar proteins are not excluded (Fig. 1A-B), the exclusion of structures containing test-set similar ligands has a consistently deleterious effect on the performance of all scoring functions, regardless of whether they make use of the RDKit descriptors. The performance of each scoring function on each of the 2007, 2013, and 2016 core sets, as well as bootstrapped 95% confidence intervals for the value of  $\rho_p$  when training on the PDBbind 2018 refined set are included in the Supporting Information (SI Fig. 9). For all scoring functions under all training and test scenarios, we reject the null hypothesis that the correlation between predicted and experimental affinity was due to random chance (permutation test  $p < 0.05$ ).

Using a RF with only RDKit descriptors consistently results in a scoring function with greater performance than a RF with the AutoDock Vina terms, and is often close to or even exceeds the performance of RF-Score, even when test-set similar ligands are excluded from the training data. This is surprising, since we would not expect to be able to predict protein-ligand binding affinity across a diverse set of protein-ligand complexes without knowing which protein the ligand was assayed with. We

discuss the possible source and interpretation of this signal next, but this observation may suggest that the redundancy between the PDBbind core set and the remainder of the database makes this particular approach to training and test set construction inappropriate for validating and benchmarking machine learning scoring functions for predicting protein-ligand binding affinity.

## Ligand-Based Features are Predictive of Mean Binding Affinity

Given the success of the RDKit descriptors alone (Fig. 1) we examined the affinity predictions for ligands that bind to proteins in multiple structures in the PDBbind database. Our RDKit RF model will produce only one value so it must be “incorrect” for many of the protein-ligand complexes. When the RDKit RF model was tested on a previously-unseen ligand, having been trained on all other data in the PDBbind 2017 general set, the score was found to be strongly correlated with the mean experimental  $pK$  of that ligand for its targets across the PDBbind 2017 general set ( $\rho_p = 0.72$ , Fig. 2). For the most common ligands in the PDBbind 2017 general set, the reported affinity values can span several orders of magnitude, so the RDKit RF model is not simply predicting a single “correct” value for a ligand that happens to have many similar affinity measurements. Fig. 3 shows the predictions of the RDKit RF model for the most common ligands in the PDBbind 2017 general set (markers) against the  $pK$  values reported by PDBbind (swarm plots). With the clear exception of biotin (BTN), the marker is often close to the centre of the swarm plot, indicating that the RDKit RF model appears to be predicting the mean affinity of the ligand even when the experimental data have a range of several  $pK$  units. We verified that many of the structures for biotin were biotin-streptavidin or biotin-avidin complexes, explaining the incredibly large experimental affinity values. Since ensemble-based methods such as RF cannot extrapolate beyond the range of values seen during training, it is unsurprising that the model cannot predict such a high average affinity when no such examples are

included in the training set.

Further, we found that when nearly-identical ligands (Tanimoto similarity  $> 0.9$ ) were excluded from the training data, this correlation is actually stronger ( $\rho_p = 0.75$ ) (Supporting Information Fig. 14). Reducing the Tanimoto similarity above which ligands are excluded from the training set results in gradually weakening correlation (Supporting Information Fig. 14), suggesting that the RDKit RF model is not simply memorizing the average affinity of highly similar training ligands.

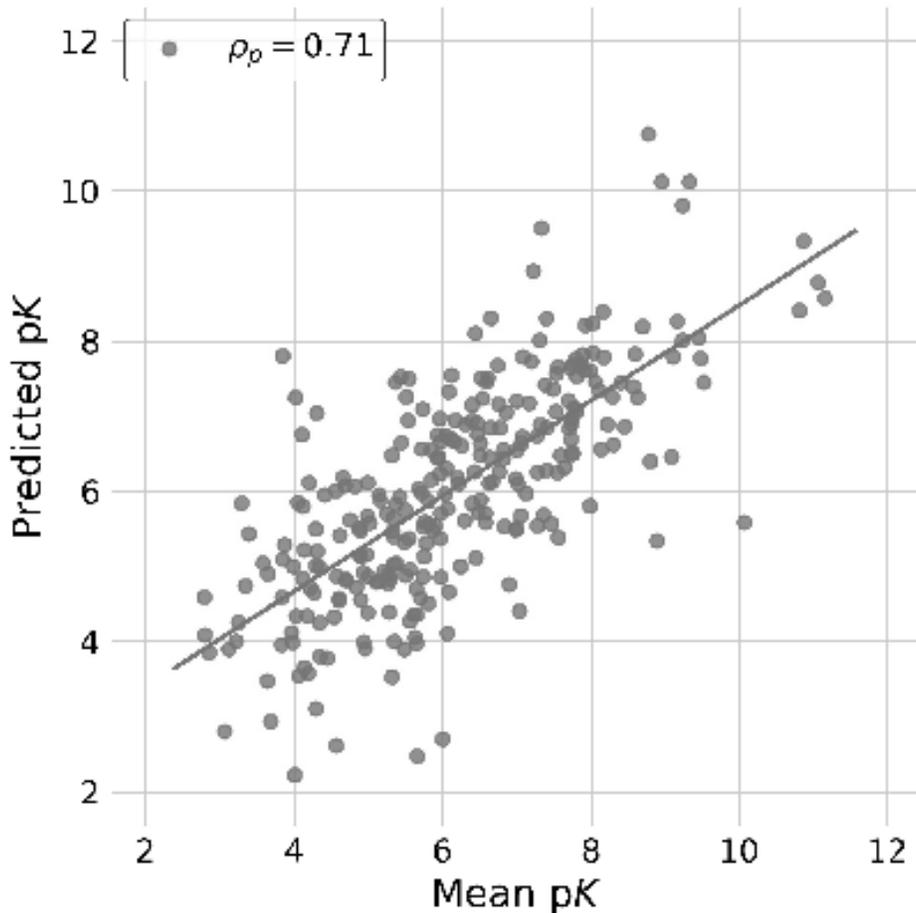


Figure 2: Correlation of ligand-based affinity prediction with mean ligand affinity. Ligands with identical chemical ID were excluded from the training data. Pearson correlation coefficient  $\rho_p = 0.72$ . Each point represents one ligand. For each ligand, a new RF regression was trained.

## Both Structure-Based and Ligand-Based Features are Important

The relative importance of the twenty highest-ranked ligand-based features for the RDKit RF model trained on the PDBbind 2018 refined set is shown in Fig. 4. The bulk properties molar refractivity (MolMR) and the logarithm of the octanol-water partition coefficient (MolLogP) are ranked highest. Molar refractivity captures the total polarizability of the molecule and log P captures its solubility; we might therefore expect these features to capture useful information about the ability of a small molecule to bind to a charged, buried active site. Both these properties are also used to characterize drugs<sup>41</sup> and log P is also used to predict bioavailability<sup>42</sup>. It is possible that the predictive power of these features can in part be attributed to systematic bias in favour of crystallising complexes featuring high-affinity engineered compounds. However, we found that there was no trivial correlation between either of these features and the  $pK$  of a compound ( $\rho_p = 0.26$  and  $\rho_p = 0.16$ , respectively, across the PDBbind 2018 general set), suggesting that their contribution to the scoring function is only in concert with other features. Perhaps easier to explain are features capturing

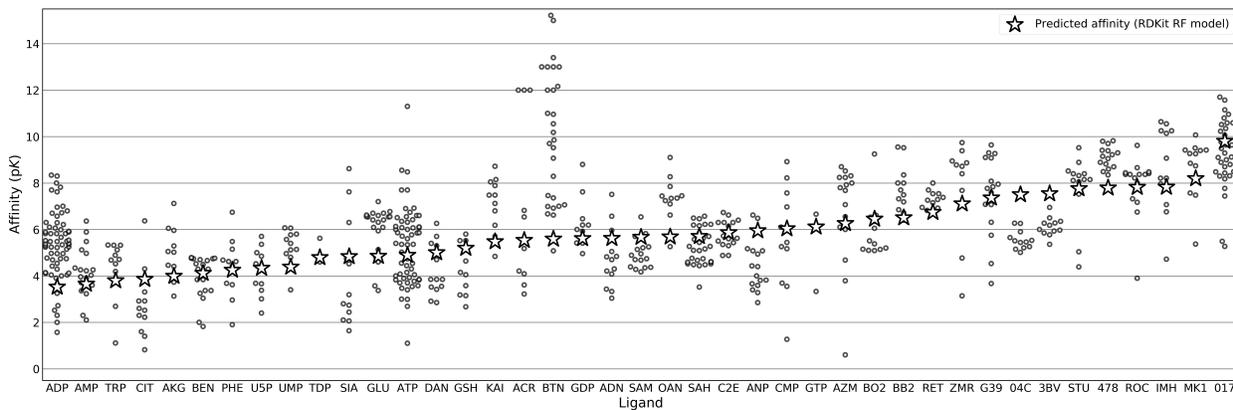


Figure 3: Ligand feature-based affinity predictions (white star) against experimental values (dots) for the most common ligands in the PDBbind 2017 general set.

molecular weight and charge (ExactMolWt, MaxAbsPartialCharge, MolWt, MinAbsPartialCharge, MaxPartialCharge) as the size and charge of the molecule will impose constraints on both its ability to fit within a binding pocket, its electrostatic com-

plementarity, and the number of interactions it has the ability to form. Similarly, topological polar surface area (TPSA) is an approximation of a molecule’s polar surface area computed using its 2D chemical graph, and may provide information about its hydrophobicity and ability to fit within a binding pocket. Van der Waals surface area contributions, captured by the PEOE\_VSA descriptors, likewise characterise the molecular surface and hence potential interactions. More complicated are the 2D descriptors capturing molecular connectivity (Chi) and graph complexity (BertzCT), whose contribution to the model might also be through capturing the shape and surface area of the molecule or some aspect of conformational entropy. We also found that when ligand-based features were combined with structure-based features in our other models, both ligand-based and structure-based features were ranked highly, and that the same ligand-based features were consistently found to be important regardless of which structure-based features were used (Supporting Information Fig. 5). This suggests that the ligand-based features are consistently capturing useful information that is not present in the structure-based features, beyond the count of rotatable bonds in the ligand.

## Conclusion

We have shown that the inclusion of a diverse set of readily-computed ligand-based features in machine learning scoring functions consistently improves their ability to rank ligands by their protein-ligand binding affinity.

Varying the composition of the training set chronologically, by restricting to data available only up until a particular year, had little effect on affinity predictions. This suggests an element of learning saturation for the targets tested with the data currently available. We showed that, in contrast, excluding proteins from the training set that are sequence-similar to those in the test set has a deleterious effect on affinity predictions and that even excluding only those proteins with identical sequence to those in the test set leads to significantly reduced scoring function performance. We also showed that

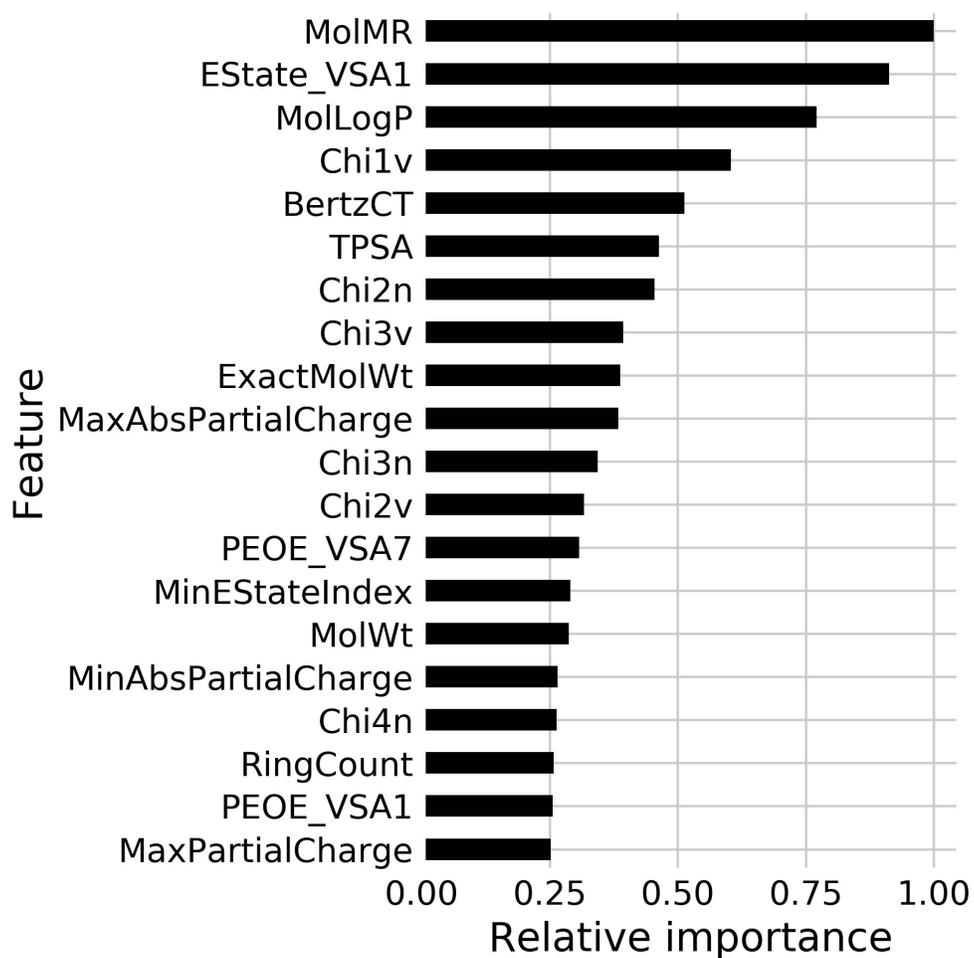


Figure 4: Relative importance of features in the RDKit RF model trained on the PDBbind 2018 refined set. A description of each feature is provided in the RDKit documentation (<https://www.rdkit.org/docs/GettingStartedInPython.html>, accessed 17/05/2019).

even when ligands with high Tanimoto similarity to those in the test set were excluded from the training set, the predictive power of the scoring functions was still increased by including ligand-based features.

Given the power of the ligand-based features, we investigated their predictive ability for ligands that bind to multiple targets and found that the predicted binding affinity of a model using only ligand-based features was strongly correlated with the mean of the experimental protein-ligand binding affinity of a ligand for its binding partners. This correlation remained strong when ligands with a Tanimoto similarity of greater than 0.9 to the test ligand were excluded from the training data. This correlation gradually weakened when progressively less similar ligands were also excluded, suggesting that while the model’s predictions are not reliant upon memorization of previously-seen highly-similar ligands, it does not extrapolate well to completely novel ligands.

Finally, we analysed the relative importance of the features of each scoring function. We found that when structure-based and ligand-based features are combined, both structure-based and ligand-based features were ranked highly, and that the same ligand-based features are ranked highly regardless of which structure-based features they were combined with. This suggests that the same information is consistently extracted from the ligand-based features and that this information is not redundant with that provided by structure-based features. Our results suggest that even under stringent validation, the addition of a diverse, quickly-computed set of ligand-based features to a scoring function yields improved predictions of binding affinity.

## Funding

This work was supported by funding from the Engineering and Physical Sciences Research Council (EPSRC) [grant number EP/G03706X/1].

## References

- (1) Gilson, M. K.; Zhou, H. X. Calculation of protein-ligand binding affinities. *Annual Review of Biophysics and Biomolecular Structure* **2007**, *36*, 21–42.
- (2) Huang, S.-Y.; Grinter, S. Z.; Zou, X. Scoring functions and their evaluation methods for protein–ligand docking: recent advances and future directions. *Physical Chemistry Chemical Physics* **2010**, *12*, 12899–12908.
- (3) Sousa, S. F.; Fernandes, P. A.; Ramos, M. J. Protein-ligand docking: current status and future challenges. *Proteins* **2006**, *65*, 15–26.
- (4) Sousa, S. F.; Ribeiro, A. J.; Coimbra, J. T.; Neves, R. P.; Martins, S. A.; Moorthy, N. S.; Fernandes, P. A.; Ramos, M. J. Protein-ligand docking in the new millennium—a retrospective of 10 years in the field. *Current Medicinal Chemistry* **2013**, *20*, 2296–314.
- (5) Ripphausen, P.; Stumpfe, D.; Bajorath, J. Analysis of structure-based virtual screening studies and characterization of identified active compounds. *Future Medicinal Chemistry* **2012**, *4*, 603–13.
- (6) Aldeghi, M.; Heifetz, A.; Bodkin, M. J.; Knapp, S.; Biggin, P. C. Accurate calculation of the absolute free energy of binding for drug molecules. *Chemical Science* **2016**, *7*, 207–218.
- (7) Perez, A.; Morrone, J. A.; Simmerling, C.; Dill, K. A. Advances in free-energy-based simulations of protein folding and ligand binding. *Current Opinion in Structural Biology* **2016**, *36*, 25–31.
- (8) Jones, G.; Willett, P.; Glen, R. C.; Leach, A. R.; Taylor, R. Development and validation of a genetic algorithm for flexible docking. *Journal of Molecular Biology* **1997**, *267*, 727 – 748.
- (9) Friesner, R. A.; Banks, J. L.; Murphy, R. B.; Halgren, T. A.; Klicic, J. J.; Mainz, D. T.; Repasky, M. P.; Knoll, E. H.; Shelley, M.; Perry, J. K.; Shaw, D. E.;

- Francis, P.; Shenkin, P. S. Glide: A New Approach for Rapid, Accurate Docking and Scoring. 1. Method and Assessment of Docking Accuracy. *Journal of Medicinal Chemistry* **2004**, *47*, 1739–1749.
- (10) Halgren, T. A.; Murphy, R. B.; Friesner, R. A.; Beard, H. S.; Frye, L. L.; Pollard, W. T.; Banks, J. L. Glide: A New Approach for Rapid, Accurate Docking and Scoring. 2. Enrichment Factors in Database Screening. *Journal of Medicinal Chemistry* **2004**, *47*, 1750–1759.
- (11) Abagyan, R.; Totrov, M.; Kuznetsov, D. ICMa new method for protein modeling and design: applications to docking and structure prediction from the distorted native conformation. *Journal of Computational Chemistry* **1994**, *15*, 488–506.
- (12) Rarey, M.; Kramer, B.; Lengauer, T.; Klebe, G. A fast flexible docking method using an incremental construction algorithm. *Journal of Molecular Biology* **1996**, *261*, 470–489.
- (13) Jain, A. N. Surflex: fully automatic flexible molecular docking using a molecular similarity-based search engine. *Journal of Medicinal Chemistry* **2003**, *46*, 499–511.
- (14) Morris, G. M.; Huey, R.; Lindstrom, W.; Sanner, M. F.; Belew, R. K.; Goodsell, D. S.; Olson, A. J. AutoDock4 and AutoDockTools4: Automated docking with selective receptor flexibility. *Journal of Computational Chemistry* **2009**, *30*, 2785–2791.
- (15) Trott, O.; Olson, A. J. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of Computational Chemistry* **2010**, *31*, 455–461.
- (16) Ravindranath, P. A.; Forli, S.; Goodsell, D. S.; Olson, A. J.; Sanner, M. F. AutoDockFR: advances in protein-ligand docking with explicitly specified binding site flexibility. *PLoS Computational Biology* **2015**, *11*, e1004586.

- (17) Li, Y.; Liu, Z.; Li, J.; Han, L.; Liu, J.; Zhao, Z.; Wang, R. Comparative assessment of scoring functions on an updated benchmark: 1. Compilation of the test set. *Journal of Chemical Information and Modeling* **2014**, *54*, 1700–16.
- (18) Li, Y.; Han, L.; Liu, Z.; Wang, R. Comparative assessment of scoring functions on an updated benchmark: 2. Evaluation methods and general results. *Journal of Chemical Information and Modeling* **2014**, *54*, 1717–36.
- (19) Ballester, P. J.; Mitchell, J. B. A machine learning approach to predicting protein–ligand binding affinity with applications to molecular docking. *Bioinformatics* **2010**, *26*, 1169–1175.
- (20) Durrant, J. D.; McCammon, J. A. NNScore 2.0: a neural-network receptor–ligand scoring function. *Journal of Chemical Information and Modeling* **2011**, *51*, 2897–2903.
- (21) Zilian, D.; Sotriffer, C. A. SFCscore RF: a random forest-based scoring function for improved affinity prediction of protein–ligand complexes. *Journal of Chemical Information and Modeling* **2013**, *53*, 1923–1933.
- (22) Li, H.; Leung, K.-S.; Wong, M.-H.; Ballester, P. J. Substituting random forest for multiple linear regression improves binding affinity prediction of scoring functions: Cyscore as a case study. *BMC Bioinformatics* **2014**, *15*, 291.
- (23) Wójcikowski, M.; Kukielka, M.; Stepniewska-Dziubinska, M. M.; Siedlecki, P. Development of a protein–ligand extended connectivity (PLEC) fingerprint and its application for binding affinity predictions. *Bioinformatics* **2018**,
- (24) Lin, H.; Sassano, M. F.; Roth, B. L.; Shoichet, B. K. A pharmacological organization of G protein–coupled receptors. *Nature Methods* **2013**, *10*, 140.
- (25) Sotriffer, C. A.; Sanschagrin, P.; Matter, H.; Klebe, G. SFCscore: scoring functions for affinity prediction of protein–ligand complexes. *Proteins: Structure, Function, and Bioinformatics* **2008**, *73*, 395–419.

- (26) van Westen, G. J.; Wegner, J. K.; Geluykens, P.; Kwanten, L.; Vereycken, I.; Peeters, A.; Ijzerman, A. P.; van Vlijmen, H. W.; Bender, A. Which compound to select in lead optimization? Prospectively validated proteochemometric models guide preclinical development. *PLoS One* **2011**, *6*, e27518.
- (27) Ain, Q. U.; Mendez-Lucio, O.; Ciriano, I. C.; Malliavin, T.; van Westen, G. J.; Bender, A. Modelling ligand selectivity of serine proteases using integrative proteochemometric approaches improves model performance and allows the multi-target dependent interpretation of features. *Integrative Biology* **2014**, *6*, 1023–33.
- (28) Paricharak, S.; Cortes-Ciriano, I.; AP, I. J.; Malliavin, T. E.; Bender, A. Proteochemometric modelling coupled to in silico target prediction: an integrated approach for the simultaneous prediction of polypharmacology and binding affinity/potency of small molecules. *Journal of Chemoinformatics* **2015**, *7*, 15.
- (29) Li, H.; Leung, K.-S.; Wong, M.-H.; Ballester, P. J. Improving AutoDock Vina using random forest: the growing accuracy of binding affinity prediction by the effective exploitation of larger data sets. *Molecular informatics* **2015**, *34*, 115–126.
- (30) Breiman, L. Random forests. *Machine learning* **2001**, *45*, 5–32.
- (31) Liu, Z.; Su, M.; Han, L.; Liu, J.; Yang, Q.; Li, Y.; Wang, R. Forging the basis for developing protein–ligand interaction scoring functions. *Accounts of Chemical Research* **2017**, *50*, 302–309.
- (32) Cheng, T.; Li, X.; Li, Y.; Liu, Z.; Wang, R. Comparative assessment of scoring functions on a diverse test set. *Journal of chemical information and modeling* **2009**, *49*, 1079–1093.
- (33) Su, M.; Yang, Q.; Du, Y.; Feng, G.; Liu, Z.; Li, Y.; Wang, R. Comparative Assessment of Scoring Functions: The CASF-2016 Update. *Journal of Chemical Information and Modeling* **2018**,

- (34) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E. The Protein Data Bank. *Nucleic Acids Research* **2000**, *28*, 235–242.
- (35) Li, H.; Leung, K.-S.; Wong, M.-H.; Ballester, P. Low-quality structural and interaction data improves binding affinity prediction via random forest. *Molecules* **2015**, *20*, 10947–10962.
- (36) Wójcikowski, M.; Zielenkiewicz, P.; Siedlecki, P. Open Drug Discovery Toolkit (ODDT): a new open-source player in the drug discovery field. *Journal of Cheminformatics* **2015**, *7*, 26.
- (37) Kramer, C.; Gedeck, P. Leave-Cluster-Out Cross-Validation Is Appropriate for Scoring Functions Derived from Diverse Protein Data Sets. *Journal of Chemical Information and Modeling* **2010**, *50*, 1961–1969.
- (38) Li, Y.; Yang, J. Structural and Sequence Similarity Makes a Significant Impact on Machine-Learning-Based Scoring Functions for ProteinLigand Interactions. *Journal of Chemical Information and Modeling* **2017**, *57*, 1007–1012.
- (39) Li, H.; Peng, J.; Leung, Y.; Leung, K.-S.; Wong, M.-H.; Lu, G.; Ballester, P. J. The Impact of Protein Structure and Sequence Similarity on the Accuracy of Machine-Learning Scoring Functions for Binding Affinity Prediction. *Biomolecules* **2018**, *8*, 12.
- (40) Pedregosa, F. et al. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **2011**, *12*, 2825–2830.
- (41) Ghose, A. K.; Viswanadhan, V. N.; Wendoloski, J. J. A Knowledge-Based Approach in Designing Combinatorial or Medicinal Chemistry Libraries for Drug Discovery. 1. A Qualitative and Quantitative Characterization of Known Drug Databases. *Journal of Combinatorial Chemistry* **1999**, *1*, 55–68, PMID: 10746014.

- (42) Lipinski, C. A. Lead-and drug-like compounds: the rule-of-five revolution. *Drug Discovery Today: Technologies* **2004**, *1*, 337–341.