

Does a More Precise Chemical Description of Protein–Ligand Complexes Lead to More Accurate Prediction of Binding Affinity?

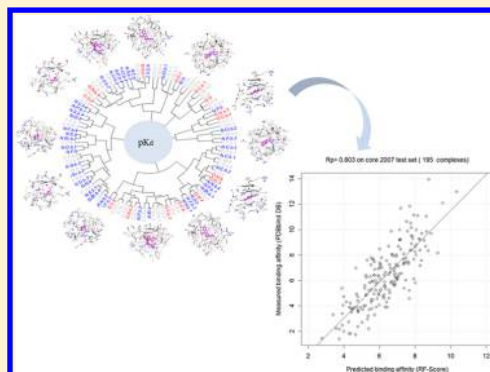
Pedro J. Ballester,^{†,*} Adrian Schreyer,[‡] and Tom L. Blundell[‡]

[†]European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton - CB10 1SD, United Kingdom

[‡]Dept. of Biochemistry, University of Cambridge, 80 Tennis Court Rd, Cambridge - CB2 1GA, United Kingdom

S Supporting Information

ABSTRACT: Predicting the binding affinities of large sets of diverse molecules against a range of macromolecular targets is an extremely challenging task. The scoring functions that attempt such computational prediction are essential for exploiting and analyzing the outputs of docking, which is in turn an important tool in problems such as structure-based drug design. Classical scoring functions assume a predetermined theory-inspired functional form for the relationship between the variables that describe an experimentally determined or modeled structure of a protein–ligand complex and its binding affinity. The inherent problem of this approach is in the difficulty of explicitly modeling the various contributions of intermolecular interactions to binding affinity. New scoring functions based on machine-learning regression models, which are able to exploit effectively much larger amounts of experimental data and circumvent the need for a predetermined functional form, have already been shown to outperform a broad range of state-of-the-art scoring functions in a widely used benchmark. Here, we investigate the impact of the chemical description of the complex on the predictive power of the resulting scoring function using a systematic battery of numerical experiments. The latter resulted in the most accurate scoring function to date on the benchmark. Strikingly, we also found that a more precise chemical description of the protein–ligand complex does not generally lead to a more accurate prediction of binding affinity. We discuss four factors that may contribute to this result: modeling assumptions, codependence of representation and regression, data restricted to the bound state, and conformational heterogeneity in data.



■ INTRODUCTION

Docking can play a key role in addressing a number of important problems such as protein–function prediction^{1,2} or drug-lead identification and optimization.^{3,4} This technique can be regarded as a two-stage process. The first is pose generation, which starts with the determination of the position, orientation, and conformation of a molecule as docked to the target's binding site. The second stage is scoring, which predicts how strongly the docked pose of such a putative ligand binds to the target. While pose generation is relatively well handled by current algorithms, the inaccuracies of current scoring functions still constitute the main barrier to achieving reliability in docking.^{5–7} Indeed, despite intensive research over more than two decades, accurate prediction of the binding affinities of large sets of diverse protein–ligand complexes remains one of the most important open problems in computational bioscience.

Three classes of scoring functions have emerged over the years: force field,^{8,9} knowledge-based,^{10–14} and empirical.^{15–19} For the sake of efficiency, scoring functions do not attempt to simulate certain physical processes that influence the process of binding. This has an impact on their ability to rank-order and selects small molecules by predicted binding affinity. Thus, two major sources of error in scoring functions arise from their

limited description of protein flexibility and the implicit treatment of solvent. Instead of scoring functions, other computational methodologies based on molecular dynamics or Monte Carlo simulations can be used to model protein flexibility and desolvation upon binding. In principle, a more accurate prediction of binding affinity than that from scoring functions is obtained in those cases amenable to these techniques.^{20,21} However, such expensive free energy calculations are not feasible for the evaluation of large numbers of protein–ligand complexes, and their application is generally limited to predicting binding affinity in series of congeneric molecules binding to a single target.²²

In addition to these two enabling simplifications, there is a third factor in scoring function development that, despite its importance, has received little attention until recently.²³ Each scoring function assumes a predetermined functional form relating the variables that describe the complex, which may also include a set of weights that are fitted to experimental or simulation data, and its predicted binding affinity. Such a relationship typically takes the form of a sum of weighted physicochemical contributions to binding in the case of

Received: February 13, 2014

Published: February 16, 2014

empirical scoring functions or a reverse Boltzmann methodology in the case of knowledge-based scoring functions. As previously discussed,²³ the inherent drawback of this approach is that those complexes not conforming to this strong modeling assumption will be predicted poorly.

As an alternative to these classical scoring functions, nonparametric machine learning can be used to capture implicitly the binding interactions that are challenging to model explicitly. By not imposing a particular functional form for the scoring function, the collective effect of intermolecular interactions in binding can be directly inferred from experimental data. The latter should lead to scoring functions with greater generality and prediction accuracy. Indeed, this unconstrained approach had to result in performance improvement given sufficient data, as it is well-known that the strong assumption of a predetermined functional form for a scoring function constitutes an additional source of error (e.g., imposing an additive form for the considered energetic contributions²⁴). On the other hand, recent experimental results have resulted in novelties in the definition of molecular interactions such as the hydrogen bond²⁵ and the hydrophobic interaction,²⁶ implying that previously proposed expressions for these energetic contributions might need to be revised accordingly.

While a few classifiers exploiting X-ray crystal structural data for discriminating between binders and nonbinders of a protein target have been presented,^{27,28} it is only recently that machine learning for nonlinear regression has been shown^{23,29} to be a particularly powerful approach to build generic scoring functions. This approach has been highlighted^{30–33} as very promising for the improvement of scoring functions. Indeed, a growing number of studies showing the benefits of machine learning scoring functions have been presented.^{23,29,34–37} However, these initial models are relatively coarse in the description of the complex, and thus the question remains as to whether the incorporation of additional chemical information relevant for binding would improve performance further.

Here, we investigate the impact of a more precise chemical characterization of the protein–ligand complex on the predictive power of the resulting scoring function. This includes the use of structural interaction fingerprints,³⁸ using atom and interaction type definitions from the CREDO structural interactomics database.³⁹ We show that the new version of RF-Score performs much better than classical scoring functions on the same test set. The RF-Score performed best when describing a complex using a 12 Å distance cutoff between atom pairs, suggesting that there is a minor contribution from long-range atom pairs. In the light of the improved performance obtained and considering the uncertainty introduced by the static nature of crystal structures, we discuss the role of interatomic distance cutoffs and binning as well as protonation states in binding affinity prediction. As a byproduct of this systematic battery of numerical experiments, the most accurate scoring function to date on a widely used pre-existing benchmark is presented. An important conclusion of our study is that a more chemically precise description of the protein–ligand complex does not generally lead to more accurate prediction of binding affinity. We discuss four factors that may contribute to this result: modeling assumptions, codependence of representation and regression, data restricted to the bound state, and conformational heterogeneity in data.

METHODS

Defining Descriptors. Each complex was described by a vector of integer-valued descriptors or features. Three description schemes were implemented: the *Element* scheme uses the combination of the element symbols of the interacting atoms to classify the interaction, e.g., C–C or N–O. The fingerprint of this scheme has a position for each pairwise combination of element symbols, and the directionality is preserved; i.e., N–O is distinct from O–N. Here, all of the heavy atoms commonly observed in PDB complexes (C, N, O, F, P, S, Cl, Br, I) are considered.

The *Sybyl* scheme uses SYBYL atom types instead of the element symbols to define the range of considered protein–ligand atom pairs. These atom types permit deconvoluting the element into hybridization state and bonding environment. For instance, instead of having a single C element atom type, the *Sybyl* scheme considers the following subtypes: C+, C1, C2, C3, Cac, and Car (a description of SYBYL atom types can be found at http://www.tripos.com/mol2/atom_types.html). The latter leads to 36 distinct C–C descriptors in the *Sybyl* scheme in contrast to a single C–C descriptor in the *Element* scheme.

The *credo* scheme uses Structural Interaction Fingerprints (SIFts)³⁸ to encode protein–ligand interactions. Here, interatomic pairs are categorized as interactions if particular geometrical and atom type constraints are satisfied. Atom types were defined through a set of SMARTS patterns that are completely customizable through a configuration file. The atom types of atoms belonging to standard amino acids as well as nonstandard binding site residues that occurred in the used test sets were precalculated because determining them “on the fly” was not feasible with the Open Babel toolkit. These atom types were stored in a separate configuration file and can therefore be easily changed by the user. The determination of standard and weak hydrogen bonds required the protonation state to be known, and there the complexes to rescore must be already protonated. Twelve different contact types are encoded in SIFT, of which four are solely distance-based. The latter are covalent bonds, van der Waals clashes and contacts, and finally proximal interactions. The other eight “feature” contact types are hydrogen bonds, weak hydrogen bonds, halogen bonds, ionic, metal complexes, aromatic, hydrophobic, and carbonyl. The definition of these including the source of the SMARTS definitions and other constraints are described in the original CREDO publication³⁹ with the following exception: the carbonyl interaction type has since then been implemented on the basis of the *ab initio* molecular-orbital calculations of Allen et al.⁴⁰ who have shown that carbonyl–carbonyl interactions can have similar strengths to those from hydrogen bonds. Appendix A3 in the Supporting Information summarizes the exact classification scheme for all interactions that was used for the SIFT descriptor.

Once the scheme is selected, descriptors are generated by counting interatomic pairs between protein binding sites and their ligand molecules. These possibly interacting atoms were assigned a specific interaction type if their distance was within a distance threshold and if a combination of possible atom type and geometry constraints was satisfied. The software to calculate these three description schemes uses the open source chemistry toolkit Open Babel⁴¹ in version 2.3.0 (through Python bindings) and the SciPy library.⁴² The molecular structures of protein targets and their ligands comprising the used data sets from PDBbind are read-in separately.

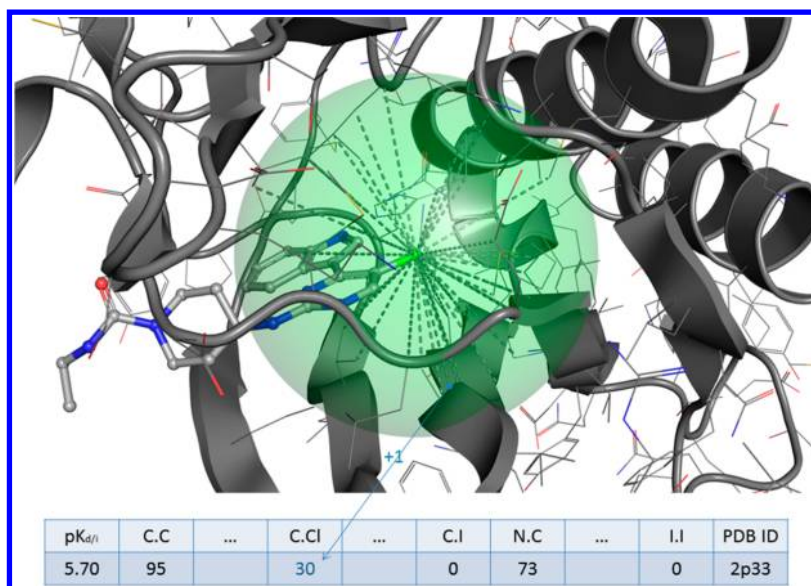


Figure 1. Sketch of the process of characterizing a protein–ligand complex (PDB: 2p33) as a set of structure-derived descriptors (C.C to I.I). The discontinuous green lines connect the ligand chlorine atom with all protein carbon atoms within the distance cutoff represented by the green sphere, with the number of these pairs giving value to the C.Cl descriptor. The rest of the descriptors are calculated in an analogous manner.

Interatomic contacts for a given distance cutoff are determined using the KDTree structure in the `scipy.spatial` module with the atomic coordinates as input (the KDTree is an efficient space-partitioning algorithm that limits the search space for interatomic contacts in order to prevent expensive all-by-all searches). The interacting atoms are then analyzed depending on the descriptor and the appropriate position on the fingerprint incremented.

This descriptor-generating software is also capable of binning the identified interatomic pairs into arbitrary distance ranges. For each normal feature on a descriptor, a number of columns equal to the number of distance bins are created. Using a distance cutoff of 6 Å and a bin size of 1 Å for example would create six bins for each feature: from 0 to 1 Å, 1–2 Å, and so on. The correct bin for each interatomic pair that has to be incremented is determined using the `numpy.digitize` function in the NumPy package (<http://numpy.scipy.org>).

The complete source code of the software that was used to generate the described results was released at <https://bitbucket.org/aschreyer/rfscore> under the MIT license.

Regression Model. RF-Score uses Random Forest⁴³ (RF) as the regression model. A RF for regression is an ensemble of P regression trees randomly generated from the same training data. In building these trees, RF determines the best split at each node of the tree from a subset of randomly selected features of size m_{try} (the recommended value of this control parameter is a third of the number of features). Note that although the selected features are generally different at each node, the same m_{try} value is applied to each node across the P trees of the forest. Here, we operate the RF with the default value $P = 500$, as its performance does not generally improve significantly beyond this threshold. The performance of each tree on predicting the Out-Of-Bag (OOB) set, here protein–ligand complexes that were not in the bootstrap sample used to train that tree, collectively provides an internal validation measure of RF. Further details about RF and its application to build RF-Score can be found in ref 23. The RF-Score software is available at <http://pedroballester.com/software>.

Training and Test Data. The PDBbind benchmark⁴⁴ has become a *de facto* standard in the validation of scoring functions. On the basis of the refined set from the 2007 release of the PDBbind database, it comprises 1300 diverse protein–ligand complexes with high quality structural and binding data (the protonation states of both proteins and ligands were already calculated by these authors⁴⁴). From this refined set, Cheng et al.⁴⁴ constructed a test set, named the core set, with 195 diverse complexes spanning more than 12 orders of magnitude in measured binding affinities. The PDBbind benchmark consists of testing the predictions of scoring functions on the core set, whereas the remaining 1105 refined set complexes are used as the training set. A discussion on the composition and suitability of this benchmark can be found in refs 23 and 29.

RESULTS

Preamble. We start this section with a concise description of the approach to predicting binding affinity using machine learning.²³ This process starts with the characterization of each protein–ligand complex as a set of intermolecular features or descriptors relevant to binding affinity prediction. The sketch in Figure 1 shows an example of how descriptors are generated from the X-ray crystal structure of a complex (PDB: 2p33). Each descriptor is given by the occurrence count of a particular protein–ligand atom pair within a predetermined distance range of each other. For example, a descriptor could be defined as the number of times that protein nitrogen and ligand oxygen are separated by less than a distance cutoff (d_{cutoff}). As in previous studies,²³ nine atom types commonly observed in PDB complexes were selected by considering atomic number only (C, N, O, F, P, S, Cl, Br, I) to give rise to 81 protein–ligand atom pairs which are considered as descriptors (this descriptor scheme is named here “element”).

Once the descriptor scheme (i.e., considered atom types, binning, and cutoff) is chosen, the next step is to select a source of curated structural and interaction data suitable for training and testing scoring functions. The PDBbind database⁴⁵ is an excellent choice for this purpose, with the additional advantage

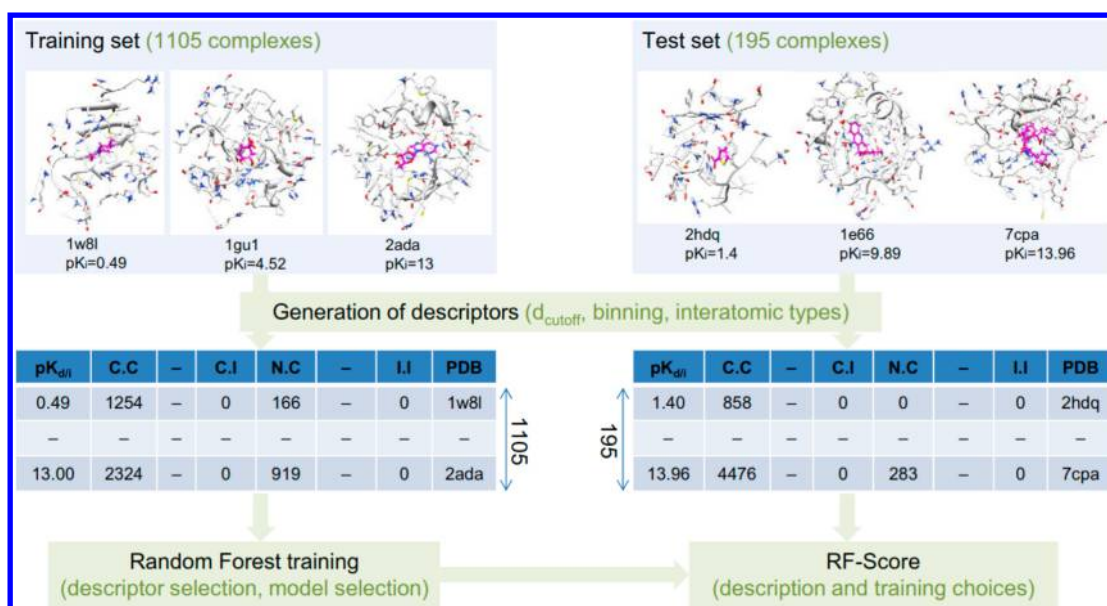


Figure 2. Training and testing RF-Score workflow. Top: descriptors are generated from two nonoverlapping data sets with 1105 and 195 complexes for training and testing, respectively. Bottom: training Random Forest to learn the nonlinear relationship between this atomic-level description of the complex and its binding affinity (pK_d or pK_i ; $pK_{d/i}$ denotes both without distinction). The resulting scoring function (RF-Score) is used to predict the binding affinity of the test set.

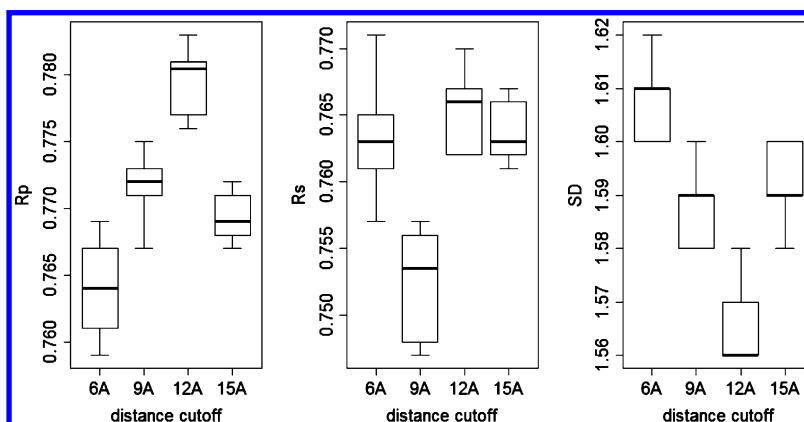


Figure 3. Test set performance of RF-Score with element descriptors for each of the four distance cutoffs (6 Å, 9 Å, 12 Å, and 15 Å). Performance is measured as the difference between observed and predicted binding affinity in the test set using three metrics: Pearson's correlation coefficient (R_p ; left plot), Spearman rank-correlation coefficient (R_s ; middle plot), and standard deviation (SD; right plot). Ten models are built from each of these four versions of the training data sets (6 Å, 9 Å, 12 Å and 15 Å), each time using a different random seed (the boxplot summarizes the performance on the test set achieved by each of the 10 models). Results showed that the best median performance, i.e., that with the highest correlations and lowest standard deviation, is obtained with the 12 Å cutoff in all three performance metrics. It is worth noting that optimizing the distance cutoff only led to a modest performance improvement (+0.017 in median R_p and $-0.05 \log K$ units in median SD).

that a large number of scoring functions have already been benchmarked on a common PDBbind test set,⁴⁴ which permits comparing new developments against the state of the art. Moreover, some scoring functions^{23,44} have not only been tested on this common data set but also calibrated on the same training set (further details can be found in the Methods section). This is important to avoid the often large bias introduced by using a different training set for each scoring function (such bias makes comparisons among scoring functions unreliable, even if compared on the same test set²⁹). Therefore, we will be focusing here on these common training and test sets.

Last, a regression model is needed to predict the binding affinity of test set complexes from the structural and interaction data in the training set. Here, we build upon RF-Score,²³ a

machine learning scoring function using RF⁴³ for regression. RF is typically tuned using a single control parameter (m_{try} , which controls the number of features that are considered for the split at each tree node) and may be subjected to a feature selection strategy intended to remove descriptors with low information content as a way to improve performance (that is, in addition to the common practice of removing all those descriptors that have zero values across training complexes). As usual, predictive performance is measured as the difference between predicted and measured binding affinity across test set complexes. Figure 2 summarizes the process of training and testing machine learning scoring functions. Full details on descriptors, data, and regression protocols can be found in the Methods section.

We have seen that not only can a protein–ligand complex be described in various ways, but also the process of building a

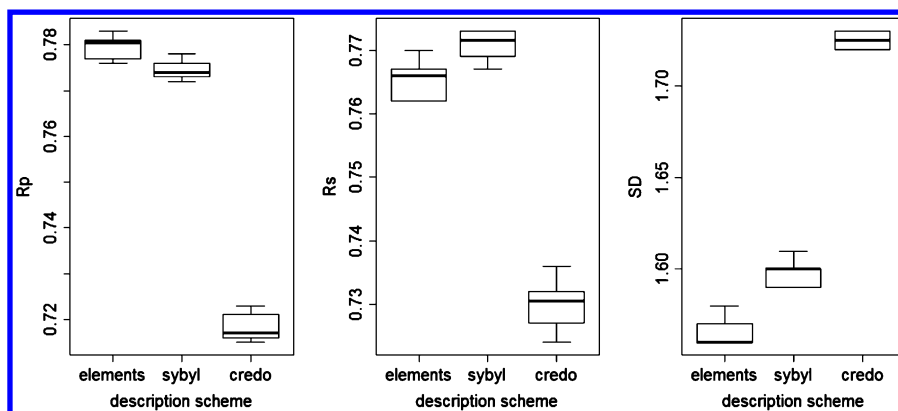


Figure 4. Test set performance of RF-Score using the optimal 12 Å interatomic distance cutoff with element, Sybyl, and credo descriptors. Interestingly, the model based on Credo descriptors obtained much lower performance than that using Sybyl and element descriptors. Element descriptors led to a small improvement over Sybyl descriptors. These results hint at a trade-off between the predictability and interpretability of the model, which we will discuss later in this paper.

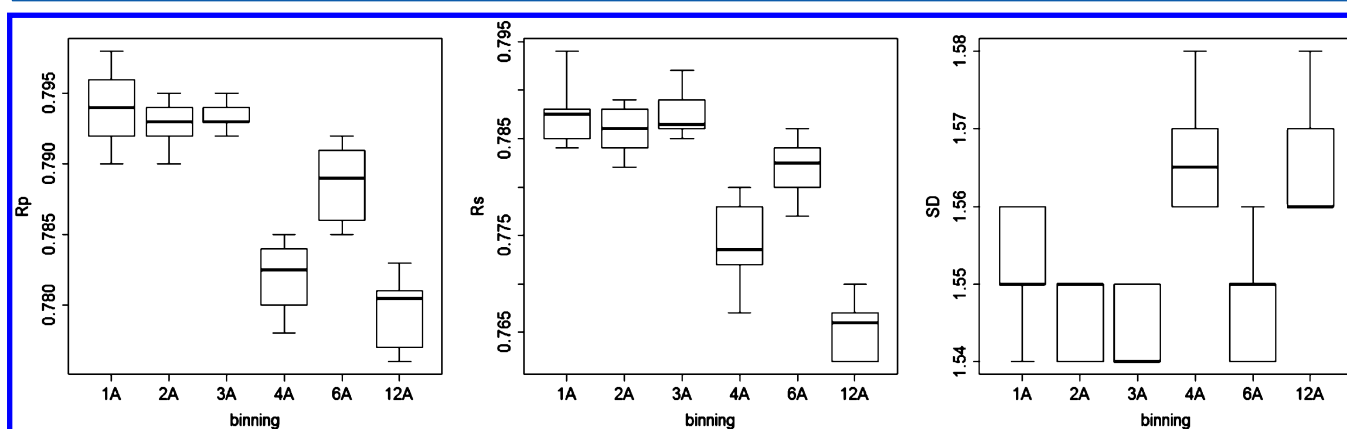


Figure 5. Test set performance of RF-Score with element descriptors and 12 Å interatomic distance cutoff using six different bin sizes (1 Å, 2 Å, 3 Å, 4 Å, 6 Å, and 12 Å). The best median performance is achieved by models with lower bin sizes (1 Å, 2 Å, and 3 Å), representing a modest improvement with respect to the model based on single-binned descriptors (12 Å).

predictive model from these descriptors involves a number of choices. In a way, the overall process of training a scoring function can be regarded as a quest for finding an optimal combination of these design variables. An exhaustive evaluation of all possible combinations is impractical, as this would involve a prohibitively large number of RF training runs (2^m runs, one for each possible subset of m features, even if we fix m_{try} to its recommended value). Thus, we were forced to assume the independence of these design variables and searched for the optimal value of each variable in a sequential manner as explained in the next subsections.

Optimal Interatomic Distance Cutoff. The first question we addressed was which distance cutoff leads to the best performance on the independent test set (henceforth referred to as simply the test set). Different distance cutoffs have been previously used in the literature, some as large as 12 Å (PMF¹²) and 15 Å (Fresno⁴⁶). Here, we addressed this question empirically by generating element descriptors for four cutoffs (6 Å, 9 Å, 12 Å, and 15 Å), which gave rise to four different numerical characterizations of the training set with their corresponding test set counterparts. Thereafter, RF was calibrated on each of these training set versions and the resulting model used to predict the binding affinity of the corresponding test set complexes. As scoring function calibration is a stochastic process, a slightly different model is

obtained with a different random seed. To assess the variability in RF prediction due to this factor, we repeated the training 10 times for each cutoff, each time with a different random seed. Such assessment is needed to establish whether the improvement in prediction is due to using another distance cutoff or just comes from variability in model calibration (this procedure is more accurate than basing the analysis on a single model calibration as has been the case so far). Lastly, we considered three commonly used metrics for quantifying the difference between predicted and measured binding affinity across the test set of protein–ligand complexes: Rp (Pearson's correlation coefficient), Rs (Spearman's correlation coefficient), and SD (Standard Deviation in log $K_{d/i}$ units). Figure 3 illustrates the results of this numerical experiment.

Role of Protonation States and Bonding Neighborhood. The element descriptors used in the previous experiment constitute a coarse representation of the complex. Distinguishing between atoms of the same element in different local environments leads to a more chemically precise characterization (e.g., deconvoluting the occurrence counts of a carbon–carbon intermolecular pair into pairs that specify the hybridization state of both atoms), and thus the resulting model would in principle be expected to perform better. To test this hypothesis, we used Sybyl atom types with these characteristics. Further, one could also incorporate additional information into

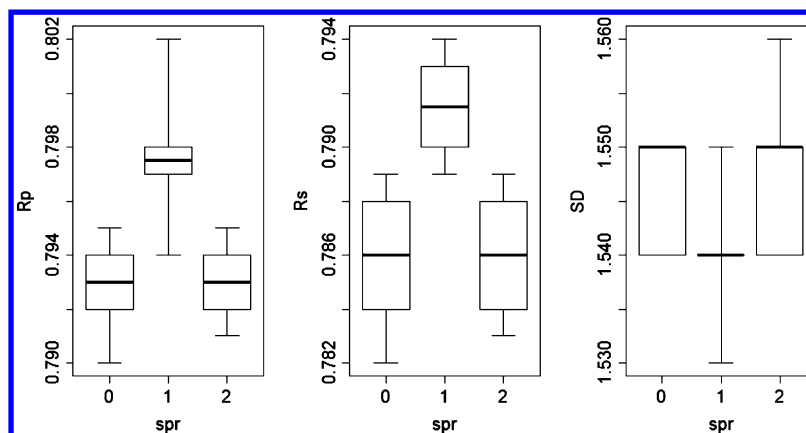


Figure 6. Test set performance of RF-Score with element descriptors, 12 Å cutoff, and 2 Å bin size using three values of the feature selection threshold (spr). Best median performance is obtained by spr = 1, which corresponds to only considering descriptors that have an average of at least one atom–atom pair in the considered distance range per training complex. The latter represents a moderate improvement with respect to the models using descriptors with spr = 0 and spr = 2.

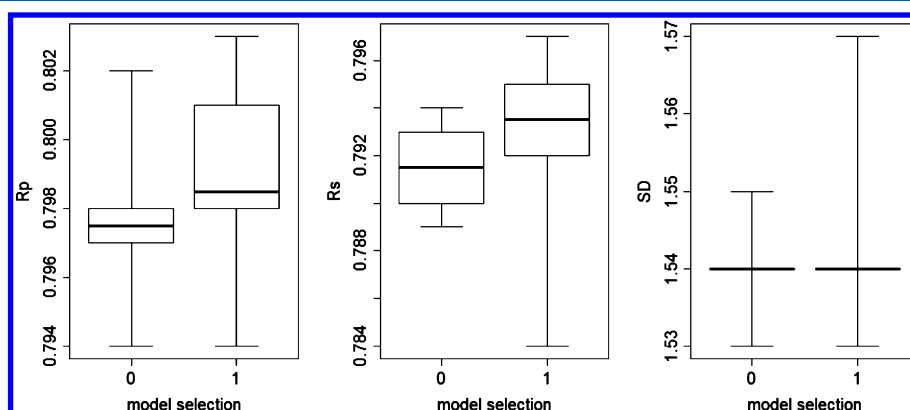


Figure 7. Test set performance of RF-Score with element descriptors, 12 Å cutoff, 2 Å bin size, and feature selection threshold (spr=1) for the recommended m_{try} (model selection=0) and m_{try} selected by OOB validation (model selection=1), which requires 123 times more training than just using the recommended setting (as many RF training runs as features were selected to describe each complex).

the descriptors such as the angle between hydrogen bond donors, acceptors, and hydrogen atoms as well as covalent and van der Waals radius. These are the Credo descriptors, which measure the abundance of a range of intermolecular interactions such as hydrogen bonds, hydrophobic interactions, or van der Waals clashes. Using the same training and test set, each description scheme gives rise to a different set of features that are used to characterize every complex. The performance of each of these three description schemes is presented in Figure 4.

Incorporating Interatomic Distance. The strength of the interatomic interactions that collectively form the noncovalent intermolecular bond depends on the separation between the interacting atoms. Therefore, it is reasonable to think that partitioning the descriptors into a number of interatomic distance bins should lead to a model with more predictivity. Consequently, we generated element descriptors with a 12 Å cutoff, i.e. using all the optimal values, for six bin sizes (a 12 Å bin size with a 12 Å cutoff simply corresponds to the case without binning, which was previously shown in the first boxplot in Figure 4 and the third boxplot in Figure 3). Figure 5 shows the results for each bin size, where the best median performance is achieved by models with lower bin sizes (1 Å, 2 Å, and 3 Å), representing a moderate improvement with respect to the model based on single-binned descriptors (12 Å).

The experiments were repeated using the same bin sizes but now with Sybyl and Credo descriptors instead of elements descriptors (the maximum cutoff for a Credo interaction type is 4.5 Å, all other atom pairs in this description scheme are labeled as “proximal”). It was observed that the performance was not as high as that with element descriptors (the best median performance for Sybyl was $R_p = 0.779$, $R_s = 0.771$, and $SD = 1.59$, whereas that for Credo was $R_p = 0.739$, $R_s = 0.742$, and $SD = 1.68$).

Feature Selection. In addition to exploring the impact of different ways to describe the complexes, we also applied basic feature selection strategies intended to remove sparse features that increased the complexity of the model without improving performance. Here, the sparsity (spr) of a descriptor is defined as the average number of occurrence counts per training complex. In previous versions of RF-Score, only features with sparsity higher than the zero threshold (spr = 0), i.e., those that are nonzero for at least one training complex, were considered. Here, we also considered two additional sparsity thresholds (spr = 1 and spr = 2). We conducted this experiment for the three best bin sizes in Figure 5 (1 Å, 2 Å, and 3 Å), which had spr = 1 as the optimal value on all three bin sizes. Figure 6 presents the results for the best bin size across the three spr values (2 Å).

Model Selection. Last, we have been using the recommended value for the RF m_{try} parameter so far. However, interval validation strategies can be used to select an optimal value for m_{try} . One of these strategies is called Out-Of-Bag (OOB) validation and essentially consists of training the model for each possible value of m_{try} and selecting the model that best performs on the internal validation set (a subset of the training set, as further explained in the Methods section). Figure 7 shows that this model selection strategy carries a small improvement in performance at the cost of much higher computational expense in model selection (one RF training run per considered m_{try} value).

Predictive Performance. This systematic battery of numerical experiments led to the new scoring function RF-Score::Elem(c12,b2)_spr1_oob (RF-Score::Elem-v2 for short). As we have seen, the descriptors come from partitioning occurrence counts of each element atom type pairs into six interatomic distance bins of 2 Å size, and the RF model is built with the 123 descriptors that are sufficiently dense ($\text{spr} = 1$) using the internally validated m_{try} value ($m_{\text{try}} = 14$ in this case). Figure 8 shows the predictive power of RF-Score::Elem-v2 on the test set.

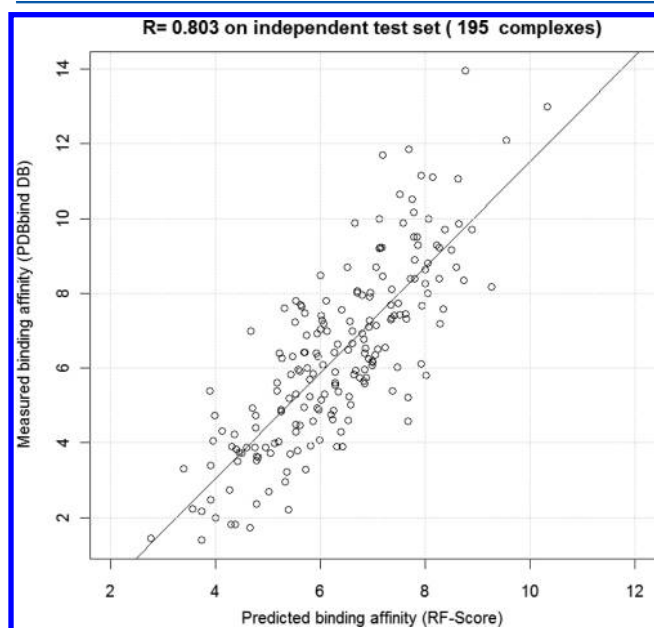


Figure 8. RF-Score::Elem-v2 predicted versus measured binding affinity on the independent test set (195 complexes). Pearson's correlation coefficient $R_p = 0.803$, Spearman's correlation coefficient $R_s = 0.797$, standard deviation $SD = 1.54 \log K_{d/i}$ units, and Root Mean Square Error $RMSE = 1.53 \log K_{d/i}$ units. This plot can be visually compared to those for the best performing scoring functions in Cheng et al.'s⁴⁴ Figure 6. Performance comparisons on the same test set are presented here in Figure 9.

In terms of efficiency, RF-Score::Elem-v2 scored all 195 protein–ligand complexes in 0.01 s (all of the computation in this study was carried out with a single processing core Intel Core i7–2920XM at 2.50 GHz with 16 GB RAM). In addition, the time to generate these features for the 195 complexes was 8 s, and hence this is the most expensive part of the calculation. Therefore, the average time to score one protein–ligand complex is about 0.04 s if the features have not been calculated before, which makes RF-Score suitable to rescore a high number of docking poses in virtual screening applications.

The predictive power of RF-Score::Elem-v2 was also compared against that of a wide selection of scoring functions on the PDBbind benchmark.⁴⁴ By using a pre-existing benchmark, the danger of constructing a benchmark complementary to the presented scoring function is avoided. It also has the advantage of ensuring that previously tested scoring functions were provided with optimal settings by their authors. Several of the scoring functions tested in the PDBbind benchmark have different versions or multiple options. However, for the sake of practicality, only the version/option of each scoring function that performs best on the PDBbind benchmark was reported by Cheng et al.⁴⁴ In addition to these 16 scoring functions, we also tested a more recent function called IMP::RankScore.⁴⁷ Figure 9 reports the performance of all scoring functions on the test set, with RF-Score::Elem-v2 obtaining the best performance with $R_p = 0.803$ (the performance of the original version of RF-Score²³ is also included). In contrast, classical scoring functions tested on the same test set obtained a lower R_p spanning from 0.216 to 0.644. This trend was also observed with the other two performance measures (R_s , SD). It is worth noting that the root-mean-square error of the free energy of binding on such a diverse test set is just 2.1 kcal/mol.

When introducing a scoring function, only the scoring function built with the random seed that provides the best performance is generally reported. We have followed here a more precise way to assess performance differences between scoring functions by comparing median performances from a set of independent trials. Moreover, in order to address the question of how significant is the reported improvement over the original version of RF-Score, we have trained and tested the original RF-Score using 10 different random seeds. Thereafter, we have repeated the process, using the same random seeds, with the new version of RF-Score. The resulting boxplots are compared in Figure 10. Lastly, we carried out a two-sample t test for each performance measure to find out that all differences are statistically significant (R_p p value = 6.0×10^{-12} , R_s p value = 1.5×10^{-11} , and SD p value = 3.8×10^{-4}).

DISCUSSION

The new version of RF-Score performs much better than classical scoring functions on the same test set. In fact, this performance gain must be actually larger in most cases since only RF-Score and X-Score, among all scoring functions represented in Figure 9, use training sets that do not overlap with the test set. Having training complexes in the test set artificially enhances the performance of a scoring function, as it is not exclusively predicting unseen complexes but merely reporting the lower training error of those overlapping complexes. Indeed, adding a third of the test set to the training set makes RF-Score::Elem-v2's R_p rise from 0.803 (no overlap between training and test sets) to 0.872 (65 overlapping complexes). Clearly, the same training and test set must be used when comparing scoring functions, but unfortunately this has not been required in recent community benchmarks.⁴⁸

Furthermore, it could be argued that there is something particular about the training/test partition selected by Cheng et al. in the PDBbind benchmark. This partition was chosen to compare RF-Score against the best scoring function in that study under exactly the same experimental conditions. An experiment to investigate this question was already carried out in the original RF-Score paper²³ (Appendix A4 in the supplementary data of that paper²³) and further discussed in

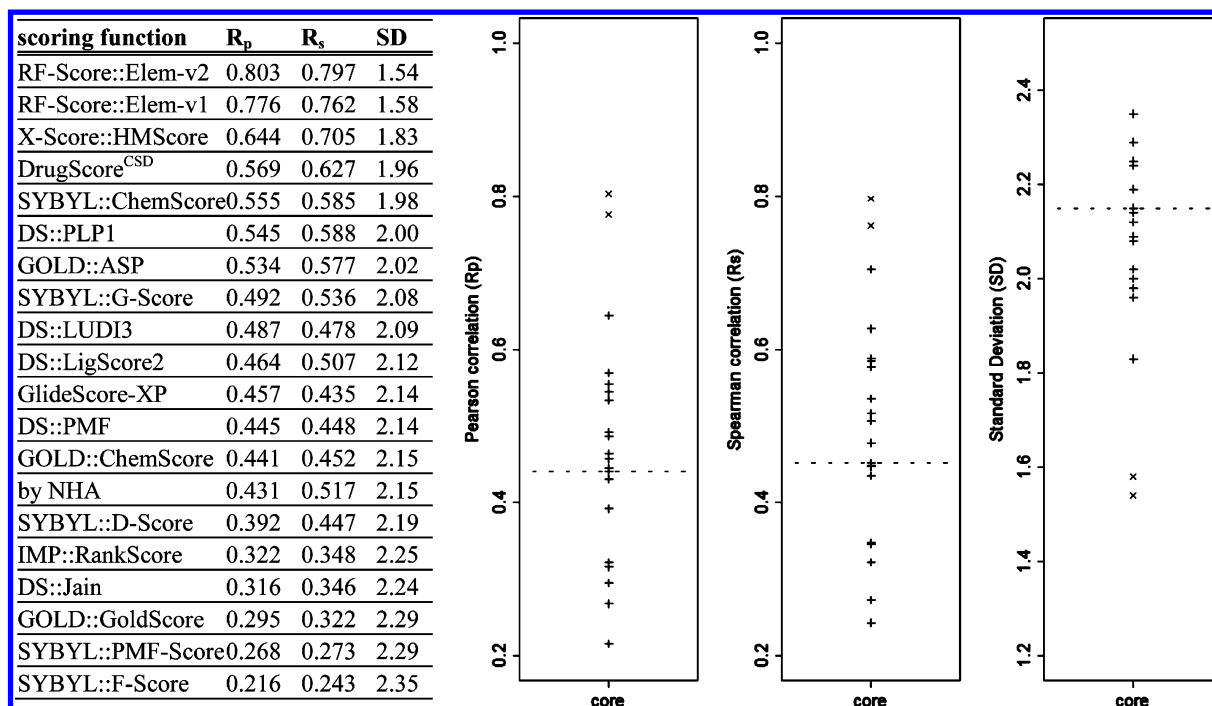


Figure 9. Performance of 18 scoring functions on the PDBbind benchmark as measured by Pearson's correlation coefficient (R_p), Spearman's correlation coefficient (R_s), and standard deviation of the difference between predicted and measured binding affinity (SD). The three plots on the right visually show the relative predictive power of RF-Score ("x" signs) against that of the other 17 scoring functions ("+" signs). NHA is the performance of a linear regression model with the number of heavy atoms of the ligand as only variable (this baseline is shown as a horizontal discontinuous line in the plots).

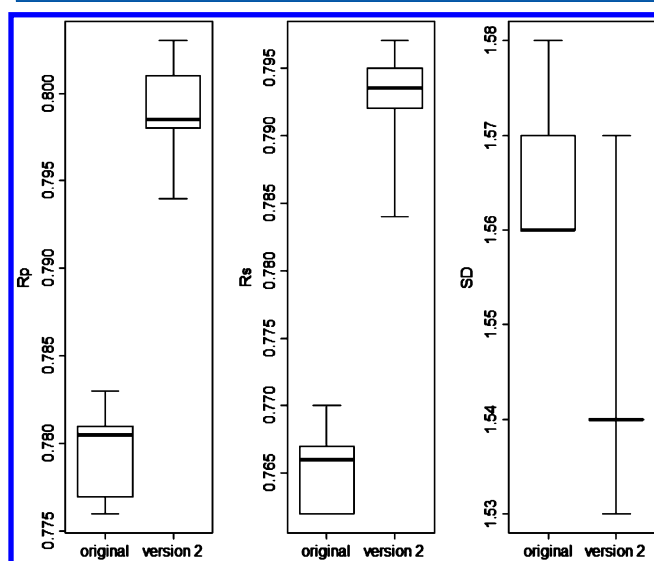


Figure 10. Performance comparison between the original version and the new version of RF-Score.

a subsequent commentary²⁹ (see Figure 1 therein). The performances of RF-Score for 25 randomly generated training/test partitions with the same sizes as the benchmark partition (1105/195) were calculated. The experiment demonstrated that there is a minor difference in RF-Score performance between the benchmark partition and the median of these 25 alternative partitions.

While our study focuses on generic scoring functions, we would like to briefly comment on how RF-Score would perform on subgroups of the test set (e.g., complexes whose proteins belong to the same family). Clearly, the better the performance

of RF-Score over another scoring function on the full test set, the more test subsets will be better predicted by RF-Score. To illustrate this, Appendix A1 presents the performance of the new version of RF-Score (RF-Score::Elem-v2; full test set RMSE = 1.54) and RF using Credo intermolecular interaction features (RF-Score::Credo; full test set RMSE = 1.72) on the four subsets resulting from partitioning the test set by binding affinity ranges. Appendix A2 presents another experiment where two small subsets of 23 complexes are generated, one containing those with the most similar ligands and the other subset with the most dissimilar ligands in terms of chemical structure. RF-Score::Elem-v2 outperforms RF-Score::Credo in all subsets but the one with the most dissimilar ligands, where RF-Score::Credo performs slightly better. These results illustrate the fact that, despite RF-Score::Elem-v2 generally performing better, there could be a few complexes (e.g., ligand-bound structures of a particular target) where other scoring functions perform better. We intend to study this issue further in the future.

On the other hand, it is noteworthy that RF-Score performed best when describing a complex using a 12 Å distance cutoff between atom pairs, a distance well beyond direct interatomic contacts. The improvement over RF-Score with a more common 6 Å cutoff is, however, modest (+0.017 in median R_p , +0.003 in median R_s , and $-0.05 \log K$ units in median SD; see Figure 3). This result suggests that there is a minor favorable contribution from atom pairs separated by a distance between 6 and 12 Å over the 1300 considered complexes. Such long-range contributions to protein–ligand binding affinity have been attributed to the electronic properties of the binding site and ligand being affected by all protein atoms⁴⁹ and also to long-range electrostatics interactions.⁵⁰ Increasing the non-covalent cutoff to 12 Å has also been found beneficial in protein

folding⁵¹ and DNA molecular dynamics⁵² simulations. It seems therefore that RF-Score is able to capture long-range effects implicitly to some extent.

The main reason why RF-Score works much better than classical scoring functions at predicting binding affinity of most complexes is due to the circumvention of the strong assumption of a predetermined functional form. All classical scoring functions consist of a sum of typically nonlinear terms with respect to selected interatomic distances, such as van der Waals terms in empirical scoring functions or particular atom–atom potentials in knowledge-based scoring functions. For instance, the Scoring Function Consortium (SFC), in a concerted effort between 10 pharmaceutical companies and academic institutions, generated an empirical scoring function (SFCscore),⁵³ which by the time of its development was clearly superior to most of the then available scoring functions. Very recently, one of the leading SFC authors has demonstrated⁵⁴ that, by using Random Forest regression instead of SFCscore's additive functional form and keeping all other modeling choices unaltered (training data, test data, and descriptors), performance rises from RMSE = 1.84 to RMSE = 1.56 (0.683 to 0.788 in the case of Rs). This is a very large improvement for a single modification in a generic model, especially taking into account that scoring functions are highly optimized due to intense work over the years in this area. Another study demonstrating that assuming an additive functional form is detrimental for the performance of empirical scoring functions is by Kinnings et al.³⁰ As force-field and knowledge-based scoring functions make the same assumption, these studies strongly suggest that a machine-learning version of other classical scoring functions will also result in significant improvement.

Another important conclusion of our study is that a more precise chemical description of the protein–ligand complex does not generally lead to more accurate prediction of binding affinity (see Figure 4). In the first study, Li et al.⁵⁵ present a scoring function tested on exactly the same test set as us, with a much larger training set that includes ours and the use of a very precise description consisting of 50 calculated descriptors falling into nine interaction categories: van der Waals, hydrogen-bonding, electrostatic, pi-system, metal–ligand bonding, desolvation effect, entropic loss effect, shape matching, and surface property matching (Table 1 on page 593 of Li et al.'s paper). Li et al. obtained SD = 1.63 and Rs = 0.779 (Table 4 at page 597 of Li et al.'s paper), whereas RF-Score originally obtained SD = 1.58 and Rs = 0.762 on the same test set. Interestingly, these authors referenced RF-Score but did not include it in the comparison or comment on why its performance was better in some performance measures despite using much simpler descriptors and less data for training.

The second independent study provides an even more direct comparison. Zilian and Sotriffer⁵⁴ used the same training set, test set, and regression model as RF-Score. The only difference between their scoring function and ours is in the 63 used descriptors, which was one of the outcomes of the industry–academia Scoring Function Consortium. These descriptors include the number of rotatable bonds in the ligand, hydrogen bonds, aromatic interactions, and polar and hydrophobic contact surfaces, among others (a full list can be found in Table 1 of page 398 of the original SFC paper⁵³). Their best scoring function achieved RMSE = 1.56 and Rs = 0.788 (Table 1 of Zilian and Sotriffer's paper), which is slightly better than the original version of RF-Score (RMSE = 1.58 and Rs = 0.762). If the modeling assumptions implied in the calculation

of chemical properties were generally accurate, we should have seen many scoring functions performing much better than RF-Score thanks to using a more precise chemical description. But we have actually seen the opposite in these two independent studies, once we compare the performances achieved by Li et al. (SD = 1.63 and Rs = 0.779) and Zilian and Sotriffer (RMSE = 1.56 and Rs = 0.788) to that of the new version of RF-Score (SD = 1.54, RMSE = 1.54 and Rs = 0.797) on the same diverse test set. The new version differs from the original version of RF-Score in that features are distance-dependent but still do not explicitly incorporate calculated protonation states.

We discuss next four convoluted factors that may contribute to this result: modeling assumptions, codependence of representation and regression, data restricted to the bound state, and conformational heterogeneity in data. The first factor is that more precise descriptors often mean making modeling assumptions that introduce additional error. For example, the protonation state of an atom needs to be estimated in order to assign its Sybyl type, but the local change in pH induced by hydrogen bond donors/acceptors in nearby residues and water molecules is usually not incorporated into scoring functions for binding affinity prediction. Similar arguments can be constructed regarding the calculation of donor-hydrogen-acceptor angles to perceive hydrogen bonds. The question remains as to how large the impact of this error is compared to that of not considering protonation states at all (the element descriptor scheme).

The second factor, often neglected, is the optimality of problem representation (description scheme) for the applied solution construction method (regression technique). From a purely chemical perspective, deconvoluting elemental atom types into their various hybridization states constitutes a more precise description of the complex. However, this scheme also results in a higher number of features and thus more sparse features. The latter are detrimental for random forest regression because as many data as possible are needed to characterize the interaction between each pair of atom types, best achieved by minimizing the number of different types defined. In practice, the definition of atom types must reflect a compromise between these two conflicting objectives, so as to ensure that the features are backed up by sufficient data to be statistically as well as chemically meaningful. This situation gives rise to a trade-off between the predictability and interpretability of the model, which is not uncommon in regression problems⁵⁶ and has also been observed here (see Figure 4).

For the sake of efficiency, scoring functions only exploit the information contained in the bound state of the complex, as represented by a crystal structure. However, binding affinity also depends on the energetic contributions from ligand and protein desolvation as well as induced fit upon binding. The third contributing factor is therefore the uncertainty about how well a particular description of the bound complex is also describing the complex just before desolvation and induced fit takes place. We speculate that descriptors whose values change less during the binding process might be more suitable for predicting binding affinity using only data about the bound state. For instance, element descriptors do not change much in general during the binding process, as a fixed cutoff will include roughly the same protein and ligand atoms just before and after binding. In contrast, protonation states will generally change significantly upon binding because of desolvation.

The last contributing factor comes from the uncertainty arising from the fact that the crystallographers deposit a single

structural model in the PDB while several different models may fit the electron density equally well.⁵⁷ The conformational heterogeneity of a complex (i.e., several bound states are possible for this complex, at least within experimental uncertainty) means that different sets of descriptors would be generated for exactly the same binding affinity. Access to the multiple structural models of a complex that are significantly different at the binding site level is likely to be helpful in deciding how to best address this issue. In particular, it would be interesting to investigate whether combining the predictions from each structure is a better strategy than simply predicting from the deposited structure.

Our finding that binding affinity can be better predicted when calculated protonation states are not explicitly incorporated into the scoring function will be certainly seen as a controversial result by most molecular modellers. We are providing next an intuitive explanation for this result. In machine learning nomenclature, the chemical description of complexes constitutes a data representation. Representations present an opportunity to incorporate domain knowledge into the problem, which in principle can help to disentangle the different explanatory factors for variation of the predicted variable (binding affinity here) and thus lead to better performance by simplifying the regression problem. However, as domain knowledge is affected by confounding factors and implemented with various degrees of efficacy, it is entirely possible to obtain better performance by incorporating less domain knowledge (i.e., introducing less noise) and hence relying more on pure inference from the data. In our problem, binding affinity is experimentally determined in solution along a trajectory in the codependent conformational spaces of the interacting molecules, whereas the structure represents a possible final state of that process in a crystallized environment. Consequently, very precise descriptors calculated from the structure are not necessarily more representative of the dynamics of binding than less precise descriptors. This means that a more precise description will not necessarily lead to a better prediction of binding affinity, as it has been proven here using Random Forest. Because the information content of a set of variables is independent of the adopted regression model, the use of an alternative regression technique should lead to the same conclusion, although this point is still to be confirmed experimentally. We cannot stress enough that we are not making any claim about the importance of protonation for pose generation in docking. Pose generation and rescoring are different problems, and so are the objectives that the corresponding scoring functions must fulfill.

In summary, we have seen that one can be easily fooled by uncertainty when investigating more accurate scoring functions. Given the unavoidable uncertainty, we believe that rigorous and systematic numerical studies are the most reliable way to make progress in predicting intermolecular binding affinity. We hope that the availability of the RF-Score software (links are provided in the Methods section) will encourage experts in the area to try to perform better on the PDBbind benchmark using alternative chemical descriptions as a way to investigate this issue further. The code permits reproducing the results obtained by RF-Score::Elem-v2 and can also be used as a template to test alternative regression techniques implemented in R. Without any modification, the RF-Score software can be employed to rescore ligands in crystal structures or docking poses. There is a range of applications in which more accurate prediction of binding affinity of a complex would be very useful,

some of them new such as replacing force-fields in molecular dynamics simulations. Other applications include structure-based virtual screening and lead optimization. In fact, applying a simpler variant of RF-Score::Elem-v1²³ to prospective virtual screening has already been found³⁷ to excel at discovering innovative inhibitors of antibacterial targets. Very recently,⁵⁸ RF-Score::Elem-v1 has been incorporated into an easy-to-set large-scale docking Web server (<http://istar.cse.cuhk.edu.hk/idock>) to carry out virtual screening of up to 17 million purchasable molecules from the ZINC database,⁵⁹ which should be upgraded soon to RF-Score::Elem-v2.

■ ASSOCIATED CONTENT

■ Supporting Information

Appendix A1 (performance on subgroups of test set by binding affinity of complexes), Appendix A2 (performance on subgroups of test set by chemical similarity of ligands), and Appendix A3 (classification scheme for all interaction types used in CREDO). This information is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*E-mail: pedro.ballester@ebi.ac.uk.

Author Contributions

P.J.B. designed the study, implemented the new version of RF-Score, and ran the numerical experiments. A.S. implemented the software to calculate descriptors. P.J.B. wrote the manuscript using the information provided by A.S. on how the descriptors were implemented. All authors discussed results and commented on the manuscript.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

We would like to thank Dr. Nidhi Tyagi (EMBL-EBI) for feedback on the manuscript. The following sources of funding are gratefully acknowledged: the Medical Research Council for a Methodology Research Fellowship (grant no. G0902106 awarded to P.J.B.). We thank the Wellcome Trust for an SDDI grant supporting A.S.

■ REFERENCES

- (1) Song, L.; Kalyanaraman, C.; Fedorov, A. A.; Fedorov, E. V.; Glasner, M. E.; Brown, S.; Imker, H. J.; Babbitt, P. C.; Almo, S. C.; Jacobson, M. P.; Gerlt, J. A. Prediction and Assignment of Function for a Divergent N-Succinyl Amino Acid Racemase. *Nat. Chem. Biol.* **2007**, *3*, 486–491.
- (2) Hermann, J. C.; Marti-Arbona, R.; Fedorov, A. A.; Fedorov, E.; Almo, S. C.; Shoichet, B. K.; Raushel, F. M. Structure-Based Activity Prediction for an Enzyme of Unknown Function. *Nature* **2007**, *448*, 775–779.
- (3) Jorgensen, W. L. Efficient Drug Lead Discovery and Optimization. *Acc. Chem. Res.* **2009**, *42*, 724–733.
- (4) Schneider, G.; Böhm, H.-J. Virtual Screening and Fast Automated Docking Methods. *Drug Discovery Today* **2002**, *7*, 64–70.
- (5) Leach, A. R.; Shoichet, B. K.; Peishoff, C. E. Prediction of Protein-Ligand Interactions. Docking and Scoring: Successes and Gaps. *J. Med. Chem.* **2006**, *49*, 5851–5855.
- (6) Moitessier, N.; Englebienne, P.; Lee, D.; Lawandi, J.; Corbeil, C. R. Towards the Development of Universal, Fast and Highly Accurate Docking/scoring Methods: A Long Way to Go. *Br. J. Pharmacol.* **2008**, *153* (Suppl.), S7–26.

- (7) Novikov, F. N.; Zeifman, A. A.; Stroganov, O. V.; Stroylov, V. S.; Kulkov, V.; Chilov, G. G. CSAR Scoring Challenge Reveals the Need for New Concepts in Estimating Protein-Ligand Binding Affinity. *J. Chem. Inf. Model.* **2011**, *51*, 2090–2096.
- (8) Huang, N.; Kalyanaraman, C.; Bernacki, K.; Jacobson, M. P. Molecular Mechanics Methods for Predicting Protein-Ligand Binding. *Phys. Chem. Chem. Phys.* **2006**, *8*, 5166–5177.
- (9) Ewing, T. J.; Makino, S.; Skillman, A. G.; Kuntz, I. D. DOCK 4.0: Search Strategies for Automated Molecular Docking of Flexible Molecule Databases. *J. Comput. Mol. Des.* **2001**, *15*, 411–428.
- (10) Mitchell, J. B. O.; Laskowski, R. A.; Alex, A.; Thornton, J. M. BLEEP- Potential of Mean Force Describing Protein-Ligand Interactions: I. Generating Potential. *J. Comput. Chem.* **1999**, *20*, 1165–1176.
- (11) Gehlhaar, D. K.; Verkhivker, G. M.; Rejto, P. A.; Sherman, C. J.; Fogel, D. B.; Fogel, L. J.; Freer, S. T. Molecular Recognition of the Inhibitor AG-1343 by HIV-1 Protease: Conformationally Flexible Docking by Evolutionary Programming. *Chem. Biol.* **1995**, *2*, 317–324.
- (12) Muegge, I.; Martin, Y. C. A General and Fast Scoring Function for Protein-Ligand Interactions: A Simplified Potential Approach. *J. Med. Chem.* **1999**, *42*, 791–804.
- (13) Mooij, W. T. M.; Verdonk, M. L. General and Targeted Statistical Potentials for Protein-Ligand Interactions. *Proteins* **2005**, *61*, 272–287.
- (14) Gohlke, H.; Hendlich, M.; Klebe, G. Knowledge-Based Scoring Function to Predict Protein-Ligand Interactions. *J. Mol. Biol.* **2000**, *295*, 337–356.
- (15) Böhm, H. J. The Development of a Simple Empirical Scoring Function to Estimate the Binding Constant for a Protein-Ligand Complex of Known Three-Dimensional Structure. *J. Comput. Mol. Des.* **1994**, *8*, 243–256.
- (16) Eldridge, M. D.; Murray, C. W.; Auton, T. R.; Paolini, G. V.; Mee, R. P. Empirical Scoring Functions: I. The Development of a Fast Empirical Scoring Function to Estimate the Binding Affinity of Ligands in Receptor Complexes. *J. Comput.-Aided Mol. Des.* **1997**, *11*, 425–445.
- (17) Wang, R.; Lai, L.; Wang, S. Further Development and Validation of Empirical Scoring Functions for Structure-Based Binding Affinity Prediction. *J. Comput. Mol. Des.* **2002**, *16*, 11–26.
- (18) 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. I. Method and Assessment of Docking Accuracy. *J. Med. Chem.* **2004**, *47*, 1739–1749.
- (19) Krammer, A.; Kirchhoff, P. D.; Jiang, X.; Venkatachalam, C. M.; Waldman, M. LigScore: A Novel Scoring Function for Predicting Binding Affinities. *J. Mol. Graphics Model.* **2005**, *23*, 395–407.
- (20) Michel, J.; Essex, J. W. Prediction of Protein–ligand Binding Affinity by Free Energy Simulations: Assumptions, Pitfalls and Expectations. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 639–658.
- (21) Mobley, D. L. Let's Get Honest about Sampling. *J. Comput. Mol. Des.* **2012**, *26*, 93–95.
- (22) Guvench, O.; MacKerell, A. D. Computational Evaluation of Protein-Small Molecule Binding. *Curr. Opin. Struct. Biol.* **2009**, *19*, 56–61.
- (23) Ballester, P. J.; Mitchell, J. B. O. A Machine Learning Approach to Predicting Protein-Ligand Binding Affinity with Applications to Molecular Docking. *Bioinformatics* **2010**, *26*, 1169–1175.
- (24) Baum, B.; Muley, L.; Smolinski, M.; Heine, A.; Hangauer, D.; Klebe, G. Non-Additivity of Functional Group Contributions in Protein-Ligand Binding: A Comprehensive Study by Crystallography and Isothermal Titration Calorimetry. *J. Mol. Biol.* **2010**, *397*, 1042–1054.
- (25) Arunan, E.; Desiraju, G. R.; Klein, R. A.; Sadlej, J.; Scheiner, S.; Alkorta, I.; Clary, D. C.; Crabtree, R. H.; Dannenberg, J. J.; Hobza, P.; Kjaergaard, H. G.; Legon, A. C.; Mennucci, B.; Nesbitt, D. J. Definition of the Hydrogen Bond (IUPAC Recommendations 2011). *Pure Appl. Chem.* **2011**, *83*, 1637–1641.
- (26) Snyder, P. W.; Mecinović, J.; Moustakas, D. T.; Thomas, S. W.; Harder, M.; Mack, E. T.; Lockett, M. R.; Héroux, A.; Sherman, W.; Whitesides, G. M. Mechanism of the Hydrophobic Effect in the Biomolecular Recognition of Arylsulfonamides by Carbonic Anhydrase. *Proc. Natl. Acad. Sci. U. S. A.* **2011**, *108*, 17889–17894.
- (27) Li, L.; Li, J.; Khanna, M.; Jo, I.; Baird, J. P.; Meroueh, S. O. Docking to Erlotinib Off-Targets Leads to Inhibitors of Lung Cancer Cell Proliferation with Suitable in Vitro Pharmacokinetics. *ACS Med. Chem. Lett.* **2010**, *1*, 229–233.
- (28) Durrant, J. D.; McCammon, J. A. NNScore: A Neural-Network-Based Scoring Function for the Characterization of Protein–Ligand Complexes. *J. Chem. Inf. Model.* **2010**, *50*, 1865–1871.
- (29) Ballester, P. J.; Mitchell, J. B. O. Comments on “Leave-Cluster-out Cross-Validation Is Appropriate for Scoring Functions Derived from Diverse Protein Data Sets”: Significance for the Validation of Scoring Functions. *J. Chem. Inf. Model.* **2011**, *51*, 1739–1741.
- (30) Kinnings, S. L.; Liu, N.; Tonge, P. J.; Jackson, R. M.; Xie, L.; Bourne, P. E. A Machine Learning-Based Method To Improve Docking Scoring Functions and Its Application to Drug Repurposing. *J. Chem. Inf. Model.* **2011**, *51*, 408–419.
- (31) Cheng, T.; Li, Q.; Zhou, Z.; Wang, Y.; Bryant, S. H. Structure-Based Virtual Screening for Drug Discovery: A Problem-Centric Review. *AAPS J.* **2012**, *14*, 133–41.
- (32) Lahti, J. L.; Tang, G. W.; Capriotti, E.; Liu, T.; Altman, R. B. Bioinformatics and Variability in Drug Response: A Protein Structural Perspective. *J. R. Soc. Interface* **2012**, *9*, 1409–37.
- (33) Sotriffer, C. *Scoring Functions for Protein–Ligand Interactions*; Gohlke, H., Ed.; Wiley-VCH Verlag GmbH & Co. KGaA: Weinheim, Germany, 2012; pp 237–263.
- (34) Das, S.; Krein, M. P.; Breneman, C. M. Binding Affinity Prediction with Property-Encoded Shape Distribution Signatures. *J. Chem. Inf. Model.* **2010**, *50*, 298–308.
- (35) Li, L.; Wang, B.; Meroueh, S. O. Support Vector Regression Scoring of Receptor-Ligand Complexes for Rank-Ordering and Virtual Screening of Chemical Libraries. *J. Chem. Inf. Model.* **2011**, *51*, 2132–2138.
- (36) Durrant, J. D.; McCammon, J. A. NNScore 2.0: A Neural-Network Receptor-Ligand Scoring Function. *J. Chem. Inf. Model.* **2011**, *51*, 2897–2903.
- (37) Ballester, P. J.; Mangold, M.; Howard, N. I.; Robinson, R. L. M. L.; Abell, C.; Blumberger, J.; Mitchell, J. B. O. Hierarchical Virtual Screening for the Discovery of New Molecular Scaffolds in Antibacterial Hit Identification. *J. R. Soc., Interface* **2012**, *9*, 3196–3207.
- (38) Deng, Z.; Chuaqui, C.; Singh, J. Structural Interaction Fingerprint (SIFt): A Novel Method for Analyzing Three-Dimensional Protein-Ligand Binding Interactions. *J. Med. Chem.* **2004**, *47*, 337–344.
- (39) Schreyer, A.; Blundell, T. CREDO: A Protein–Ligand Interaction Database for Drug Discovery. *Chem. Biol. Drug Des.* **2009**, *73*, 157–167.
- (40) Allen, F. H.; Baalham, C. A.; Lommerse, J. P. M.; Raithby, P. R. Carbonyl–Carbonyl Interactions Can Be Competitive with Hydrogen Bonds. *Acta Crystallogr., Sect. B: Struct. Sci.* **1998**, *54*, 320–329.
- (41) O'Boyle, N. M.; Morley, C.; Hutchison, G. R. Pybel: A Python Wrapper for the OpenBabel Cheminformatics Toolkit. *Chem. Cent. J.* **2008**, *2*, 1–7.
- (42) Jones, E.; Oliphant, T.; Peterson, P.; et al. SciPy: Open source scientific tools for Python. <http://www.scipy.org/>.
- (43) Breiman, L. Random Forests. *Mach. Learn.* **2001**, *45*, 5–32.
- (44) Cheng, T.; Li, X.; Li, Y.; Liu, Z.; Wang, R. Comparative Assessment of Scoring Functions on a Diverse Test Set. *J. Chem. Inf. Model.* **2009**, *49*, 1079–1093.
- (45) Wang, R.; Fang, X.; Lu, Y.; Wang, S. The PDBbind Database: Collection of Binding Affinities for Protein-Ligand Complexes with Known Three-Dimensional Structures. *J. Med. Chem.* **2004**, *47*, 2977–2980.
- (46) Rognan, D.; Lauemoller, S. L.; Holm, A.; Buus, S.; Tschinke, V. Predicting Binding Affinities of Protein Ligands from Three-Dimen-

sional Models: Application to Peptide Binding to Class I Major Histocompatibility Proteins. *J. Med. Chem.* **1999**, *42*, 4650–4658.

(47) Fan, H.; Schneidman-Duhovny, D.; Irwin, J. J.; Dong, G.; Shoichet, B. K.; Sali, A. Statistical Potential for Modeling and Ranking of Protein–Ligand Interactions. *J. Chem. Inf. Model.* **2011**, *51*, 3078–3092.

(48) Smith, R. D.; Dunbar, J. B.; Ung, P. M.-U.; Esposito, E. X.; Yang, C.-Y.; Wang, S.; Carlson, H. A. CSAR Benchmark Exercise of 2010: Combined Evaluation Across All Submitted Scoring Functions. *J. Chem. Inf. Model.* **2011**, *51*, 2115–2131.

(49) Hayik, S. A.; Dunbrack, R.; Merz, K. M. A Mixed QM/MM Scoring Function to Predict Protein–Ligand Binding Affinity. *J. Chem. Theory Comput.* **2010**, *6*, 3079–3091.

(50) Caravella, J. A.; Carbeck, J. D.; Duffy, D. C.; Whitesides, G. M.; Tidor, B. Long-Range Electrostatic Contributions to Protein–Ligand Binding Estimated Using Protein Charge Ladders, Affinity Capillary Electrophoresis, and Continuum Electrostatic Theory. *J. Am. Chem. Soc.* **1999**, *121*, 4340–4347.

(51) Piana, S.; Lindorff-Larsen, K.; Dirks, R. M.; Salmon, J. K.; Dror, R. O.; Shaw, D. E. Evaluating the Effects of Cutoffs and Treatment of Long-Range Electrostatics in Protein Folding Simulations. *PLoS One* **2012**, *7*, e39918.

(52) Norberg, J.; Nilsson, L. On the Truncation of Long-Range Electrostatic Interactions in DNA. *Biophys. J.* **2000**, *79*, 1537–1553.

(53) Sotriffer, C. A.; Sanschagrin, P.; Matter, H.; Klebe, G. SFCscore: Scoring Functions for Affinity Prediction of Protein–Ligand Complexes. *Proteins* **2008**, *73*, 395–419.

(54) Zilian, D.; Sotriffer, C. A. SFCscore(RF): A Random Forest-Based Scoring Function for Improved Affinity Prediction of Protein–Ligand Complexes. *J. Chem. Inf. Model.* **2013**, *53*, 1923–1933.

(55) Li, G.-B.; Yang, L.-L.; Wang, W.-J.; Li, L.-L.; Yang, S.-Y. ID-Score: A New Empirical Scoring Function Based on a Comprehensive Set of Descriptors Related to Protein–Ligand Interactions. *J. Chem. Inf. Model.* **2013**, *53*, 592–600.

(56) Sukumar, N.; Das, S. Current Trends in Virtual High Throughput Screening Using Ligand-Based and Structure-Based Methods. *Comb. Chem. High Throughput Screening* **2011**, *14*, 872–888.

(57) Furnham, N.; Blundell, T. L.; DePristo, M. A.; Terwilliger, T. C. Is One Solution Good Enough? *Nat. Struct. Mol. Biol.* **2006**, *13*, 184–185.

(58) Li, H.; Leung, K.-S.; Ballester, P. J.; Wong, M.-H. Istar: A Web Platform for Large-Scale Protein–Ligand Docking. *PLoS One* **2014**, *9*, e85678.

(59) Irwin, J. J.; Sterling, T.; Mysinger, M. M.; Bolstad, E. S.; Coleman, R. G. ZINC: A Free Tool to Discover Chemistry for Biology. *J. Chem. Inf. Model.* **2012**, *52*, 1757–1768.