

Improving Scoring-Docking-Screening Powers of Protein–Ligand Scoring Functions using Random Forest

Cheng Wang^[a] and Yingkai Zhang^{*,[a,b]}

The development of new protein–ligand scoring functions using machine learning algorithms, such as random forest, has been of significant interest. By efficiently utilizing expanded feature sets and a large set of experimental data, random forest based scoring functions (RFbScore) can achieve better correlations to experimental protein–ligand binding data with known crystal structures; however, more extensive tests indicate that such enhancement in scoring power comes with significant under-performance in docking and screening power tests compared to traditional scoring functions. In this work, to improve scoring-docking-screening powers of protein–ligand docking functions

simultaneously, we have introduced a $\Delta_{\text{vina}}\text{RF}$ parameterization and feature selection framework based on random forest. Our developed scoring function $\Delta_{\text{vina}}\text{RF}_{20}$, which employs 20 descriptors in addition to the AutoDock Vina score, can achieve superior performance in all power tests of both CASF-2013 and CASF-2007 benchmarks compared to classical scoring functions. The $\Delta_{\text{vina}}\text{RF}_{20}$ scoring function and its code are freely available on the web at: <https://www.nyu.edu/projects/yzhang/DeltaVina>. © 2016 Wiley Periodicals, Inc.

DOI: 10.1002/jcc.24667

Introduction

Protein–ligand docking is a computational approach that attempts to predict the binding mode between a protein receptor and a small molecule ligand as well as their binding affinity. It plays an increasingly important role in structure-based drug design as well as in functional studies of proteins. The most critical component of docking is the scoring function, which is needed to determine binding site and binding mode of a ligand on a protein,^[1] to screen virtual small-molecule libraries to identify potential leads for further inhibitor development,^[2–8] and to explicitly estimate the binding affinity between a protein and a ligand given their complex structure. Correspondingly, to assess performance of scoring functions for these different important tasks, several key metrics have been developed and adopted, including: (i) a docking power test, which evaluates the ability of the scoring function to identify the native binding site and binding mode among a set of computer generated decoys; (ii) a screening power test to evaluate the ability of the scoring function to identify a true binder for a given target from a pool of random molecules; and (iii) a scoring power test, which assesses the linear correlation between predicted and experimental binding affinities.^[9–13] Extensive retrospective and comparative studies^[10,11,14–18] indicate that some widely used scoring functions, such as GlideScore,^[19–21] can perform relatively well in docking and screening power tests, but most perform less satisfactorily in the scoring power test. Thus, the accuracy of scoring functions remains a central limitation of protein–ligand docking.

In the last few years, machine learning approaches have proven useful for many technologies in modern society, such as computer vision and natural language processing.^[22–25] In the field of biomolecular modeling, there has been significant interest to develop new protein–ligand scoring functions using

state-of-the-art machine learning methods,^[26–40] such as the random forest (RF) algorithm. By efficiently utilizing expanded feature sets and a large set of experimental data, random forest based scoring functions (RFbScore)^[26–32] have achieved significantly better correlations with experimental protein–ligand binding data for known crystal structures; however, more extensive testing indicates that this enhancement in scoring power comes with significant under-performance in docking and screening power tests compared to traditional scoring functions.^[41,42]

Random forest is an ensemble learning method based on the aggregation of numerous decision trees.^[43,44] In RFbScore, every regression tree is a nonparametric predictive model to relate structural features to binding affinities, the predicted values of which are bounded by the learning set. Thus, random forest can do interpolation but not extrapolation.^[45] Without a predetermined function form, random forest has the ability to learn complicated interactions directly from a large set of experimental data based on numerous input features. Up to now, almost all published RFbScores that predict binding affinity have used experimental protein–ligand binding data with known crystal structures as the training set alone. Thus, in retrospect, it is not surprising that RFbScores can achieve success in scoring power tests, which mostly rely on interpolation—that is, to estimate binding affinities given experimentally determined structures.^[29] Conversely, numerous

[a] C. Wang, Y. Zhang

Department of Chemistry, New York University, New York, New York 10003

E-mail: yingkai.zhang@nyu.edu

[b] Y. Zhang

NYU-ECNU Center for Computational Chemistry at NYU Shanghai, Shanghai 200062, China

Contract grant sponsor: NIH; Contract grant number: R01-GM079223

© 2016 Wiley Periodicals, Inc.

tasks in docking and screening tests depend on extrapolation—that is, to estimate binding affinities for computationally generated structures which should have weaker binding affinities. Thus, it is understandable that RFbScores would falter when applied to such decoy structures, leading to significant underperformance of RFbScores in docking and screening power tests.^[41,42]

From the above discussion, we can see that the inferior performance of RFbScores in docking tests may reflect two problems: (1) Random forest is restricted to do interpolation; (2) Use of experimental protein–ligand binding data with known crystal structures as the training set alone limits the applicability of RFbScores. In this work, to improve scoring-docking-screening powers of scoring functions simultaneously, we employ a two-pronged strategy:

One is to expand the training set. Besides enlarging the experimental dataset to include crystal structures with weak binding affinities, we have added a similar amount of computationally generated structures (decoy data) into the training set. The idea of including decoy data in the training set has been previously employed in the development of several other scoring functions.^[34,37]

The other is to employ a $\Delta_{\text{vina}}\text{RF}$ approach, in which random forest is employed to parameterize corrections to the AutoDock Vina scoring function. This is partly inspired by the recent development of the Δ -machine learning approach to predict enthalpies of organic molecules.^[46] AutoDock Vina is one of most widely used open-source docking programs, which has been successfully employed in numerous docking and screening tasks.^[47] Our $\Delta_{\text{vina}}\text{RF}$ parameterization framework aims to combine the excellent docking power of the Vina docking function with the strength of random forest for improving scoring accuracy.

Furthermore, by employing random forest for feature selection and introducing a pharmacophore-based solvent-accessible surface area (SASA) feature set, our developed $\Delta_{\text{vina}}\text{RF}_{20}$, which employs 20 descriptors in addition to the AutoDock Vina score, can achieve superior performance compared to traditional scoring functions in all tests of both CASF-2013 and CASF-2007 benchmarks, including scoring, ranking, docking and screening power tests. It should be noted that our new scoring function $\Delta_{\text{vina}}\text{RF}_{20}$ has not been incorporated into AutoDock Vina for ligand sampling, and currently it can only be used for post-scoring.

Methods

Training set of protein–ligand complexes

The training set of protein–ligand complexes for this work consists of two subsets: one is an experimental subset, which includes 3336 crystal complex structures with experimentally measured binding affinities; the other is a decoy subset, which includes 3322 computer generated decoy structures with computationally estimated binding affinities (Supporting Information Table S1). These are obtained, respectively, from the PDBbind database,^[12,48–50] which is a collection of protein–

ligand complex PDBs with experimental binding affinities, and the CSAR decoy set, which is a collection of computer generated binding poses as well as native poses for structures in the CSAR-NRC HiQ benchmark release.^[51,52] Any structure in the CASF-2007 or CASF-2013 benchmark sets,^[12,13] which will be used for the test set, is excluded from the training set.

The experimental subset consists of data from three sources: the PDBbind refined set (v2014), native poses in the CSAR decoy dataset, and weak-binding protein–ligand crystal structures (pK_d between 0.4 to 3) in PDBbind v2014 general set. Any entry in the CASF-2007 or CASF-2013 benchmark set is excluded.

The decoy subset contains decoy data for 302 protein–ligand complexes in the CSAR decoy set, excluding 41 complexes that are also in the CASF-2007 or CASF-2013 benchmark set. For each native pose, 11 decoys are selected from up to 500 decoys in the original CSAR decoy set based on the rank of AutoDock Vina score at 0%, 10%, 20%, ..., 90%, 100%, respectively. Thus, the decoy subset has a similar number of data entries as in the experimental subset.

The binding affinity for each complex in the training set is denoted as $pK_d(\text{train})$. For each entry in the experimental subset, $pK_d(\text{train})$ is the experimental binding affinity $pK_d(\text{exp})$. The binding affinity for each entry in the decoy subset should not be larger than the experimental binding affinity for the corresponding native pose. For each decoy, first we calculate $pK_d(\text{Vina})$ based on the AutoDock Vina score: if the calculated $pK_d(\text{Vina})$ of a decoy structure is less than $pK_d(\text{exp})$ of the corresponding native pose, we assign $pK_d(\text{Vina})$ as $pK_d(\text{train})$ for this decoy structure; otherwise $pK_d(\text{train})$ of the decoy structure is assumed to be at the upper limit, which is the $pK_d(\text{exp})$ of the corresponding native pose.

The $\Delta_{\text{vina}}\text{RF}$ parameterization approach

Our main idea is to employ random forest to parameterize corrections to the AutoDock Vina scoring function, and thus to take advantage of both the excellent docking power of the Vina docking function and the strength of random forest in improving scoring accuracy. The Vina scoring function consists of six components: two Gaussian steric terms, one repulsion term, one hydrogen bonding (HB) term, and one torsion count factor, and has been parameterized to improve both binding pose and affinity prediction.^[47] The original score calculated by the AutoDock Vina program is in the unit of kcal/mol, and can be converted into pK_d unit with the following formula: $pK_d(\text{Vina}) = -0.73349 E(\text{Vina})$. Thus, our overall $\Delta_{\text{vina}}\text{RF}$ scoring function can be cast into the following form:

$$pK_d(\Delta_{\text{vina}}\text{RF}) = pK_d(\text{Vina}) + \Delta pK_d(\text{RF})$$

where $\Delta pK_d(\text{RF})$ is the correction term trained by the random forest (RF) algorithm using $\Delta pK_d(\text{train})$, that is, $pK_d(\text{train}) - pK_d(\text{Vina})$.

Given a learning set $L = \{(X^{(1)}, y^{(1)}), \dots, (X^{(N)}, y^{(N)})\}$, which contains N pairs of input feature vectors $X = (x_1, x_2, \dots, x_p)$ and output values y , each regression tree in a random forest model

can be grown as follows: (1). Sample the learning set. Prior to growing a decision tree k , a bootstrap learning subset L_k^* is drawn at random from L with replacement, and the left-out data ($L - L_k^*$) constitutes the (OOB) out-of-bag subset OOB_k ; (2). Grow a single decision tree T_k . Based on the bootstrap learning subset L_k^* , T_k is constructed by recursively splitting each terminal node of the tree into two child nodes until the minimum node size is reached. For each splitting, it picks the best feature from a pool of m_{try} features. The m_{try} features are randomly selected from all p features. (3). The prediction error of the T_k is estimated using the out-of-bag subset OOB_k . After repeating steps 1-3 to grow M regression trees, the collection of all regression trees (T_k , $k = 1, \dots, M$) is considered as a predictive RF model. To make a prediction with a new input feature vector $X^{(new)}$, its predicted value is the average of predictions from all trees:^[44]

$$y^{(new)} = \frac{1}{M} \sum_{k=1}^M T_k(X^{(new)})$$

In our development of $\Delta_{vina}RF_{20}$ parameterization, the learning set L is derived from our training set that has been described above: N is 6658; the output value y is $\Delta pK_d(\text{train})$, that is, $pK_d(\text{train}) - pK_d(\text{Vina})$; the input feature vector has $p = 20$ features, which are calculated based on the corresponding protein–ligand structure in the training set. The randomForest package in R is used to build random forest models.^[53] The final $\Delta_{vina}RF_{20}$ model is built by using $M = 500$ regression trees with $m_{try} = 4$, selected based on the OOB performance of the learning set.

The twenty features in $\Delta_{vina}RF_{20}$ are listed in Table 1. There are 10 terms from the AutoDock Vina source code and 10 terms related to buried solvent-accessible surface area (bSASA). In the AutoDock Vina source code,^[47] there are a total of 58 terms implemented as listed in Supporting Information Table S2, among which 6 terms have been selected for the AutoDock Vina scoring function, including: two Gaussian steric terms, one repulsion term, one hydrogen bonding (HB) term, and one torsion count factor. These 58 Vina-implemented terms have been explored in the development of smina and user-specified custom scoring functions with linear regression.^[54] During our development of $\Delta_{vina}RF_{20}$, we first ranked 58 Vina-implemented terms based on the permutation variable importance indices %IncMSE, which is OOB mean square error (MSE) increase as a result of feature i being permuted (values are randomly shuffled). For a given feature i , it is calculated by

$$\%IncMSE_i = \frac{MSE_i^{OOB} - MSE^{OOB}}{MSE^{OOB}} \times 100\%$$

where MSE^{OOB} is the OOB MSE and MSE_i^{OOB} is the OOB MSE when feature i is permuted. More important features have higher %IncMSE values. Then we have employed a backward feature selection approach, in which the least important features are removed one by one to build random forest models, to choose the least number of features with a comparable top

performance. The selected 10 Vina-implemented terms in Table 1 include five polar interaction terms and five ligand-dependent terms.

The bSASA terms are calculated using atomic SASA changes between the unbound and bound structures: for an atom i , $bSASA_i = SASA_{i, \text{unbound}} - SASA_{i, \text{complex}}$, where atomic SASAs are calculated by the MSMS program using a probe radius of 1.0 Å.^[56] As shown in Supporting Information Table S3, nine pharmacophore types are defined for the atoms in the protein and ligand based on SYBYL atom types and neighboring atoms as in DOCK.^[57] The SYBYL atom types^[58] are converted by Pybel from the structures with hydrogen atoms added.^[59] Thus in Table 1, there are nine pharmacophore-based bSASA terms and 1 total bSASA term.

Testing set and evaluation methods

Both CASF-2013 and CASF-2007 benchmark sets^[12,13] are used as testing sets so that the results can be directly compared with other docking functions. Both datasets consist of 195 protein–ligand complexes selected from a refined dataset in their respective year's PDBbind database.^[12,48] Scoring, ranking and docking powers have been evaluated for 16 scoring functions (in Supporting Information Table S4) for the CASF-2007 benchmark set^[12], while scoring, ranking, docking and screening power tests have been carried out for 20 scoring functions (in Supporting Information Table S5) for the CASF-2013 benchmark set.^[13] In our current work, all power tests for AutoDock Vina and $\Delta_{vina}RF_{20}$ are carried out in the same way as those described in comparative assessment articles for CASF-2007 and CASF-2013,^[13] which are summarized below.

Scoring Power. The scoring power test evaluates the linear correlation between predicted binding affinity and experimental binding affinity. It is evaluated by the Pearson's correlation coefficient (R) between predicted binding affinity and experimental binding affinity and the standard deviation (SD) in regression:

$$R = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}, \quad SD = \sqrt{\frac{\sum [y_i - (a + bx_i)]^2}{N - 1}}$$

where x_i is the predicted binding affinity for i th complex; y_i is the experimental binding affinity for i th complex; a and b are the intercept and the slope of linear regression between experimental binding affinity and predicted binding affinity.

Ranking Power. The ranking power test assesses the ability of a scoring function to correctly rank the known ligands of the same target protein based on their predicted binding affinity given the poses from the crystal structures. For each benchmark, there are 65 target proteins and three known ligands for each protein. Two levels of success, namely high-level and low-level, are evaluated in CASF-2013. For the high-level, the three ligands for target protein should be ranked by predicted score as the best > the median > the poorest, while the low-level only needs to pick the best one out of three. The success rate is calculated by the number of the correctly ranked

Table 1. 20 Features in $\Delta_{\text{vina}}\text{RF}_{20}$.

No.	Feature description
<u>AutoDock Vina interaction terms^[a]</u>	
1	$\text{Non_hydrophobic}(a_1, a_2, d) = \begin{cases} 0, & a_1 \text{ or } a_2 \text{ is hydrophobic} \\ 1, & d_{\text{diff}}(a_1, a_2) < 0.5 \\ 0, & d_{\text{diff}}(a_1, a_2) \geq 1.5 \\ 1.5 - d_{\text{diff}}(a_1, a_2), & \text{otherwise} \end{cases}$
2	$\text{Hydrogen_bond}(a_1, a_2, d) = \begin{cases} 0, & (a_1, a_2) \text{ do not form hydrogen bond} \\ 1, & d_{\text{diff}}(a_1, a_2) < -0.7 \\ 0, & d_{\text{diff}}(a_1, a_2) \geq 0.4 \\ \frac{d_{\text{diff}}(a_1, a_2) - 0.4}{-1.1}, & \text{otherwise} \end{cases}$
3	$\text{Solvation}(a_1, a_2, d) = [(ASP_{a_1} + QASP \times q_{a_1})V_{a_2} + (ASP_{a_2} + QASP \times q_{a_2})V_{a_1}]e^{-\left(\frac{d}{2}\right)^2}$
4-5	$\text{Electrostatic}(a_1, a_2, d) = \frac{q_{a_1} \times q_{a_2}}{d^x}, x = 1 \text{ or } 2$
<u>AutoDock Vina ligand dependent terms</u>	
6	Number of heavy atoms
7	Number of hydrophobic atoms
8	Number of torsion
9	Number of rotors
10	Ligand length
<u>bSASA features^[b]</u>	
11	Positive
12	Negative
13	Donor-acceptor
14	Donor
15	Acceptor
16	Aromatic
17	Hydrophobic
18	Polar
19	Halogen
20	Total bSASA
<p>[a] Interaction terms are from the AutoDock Vina source code.^[54] d is the distance between two atoms, a_1 and a_2. d_{diff} is the surface distance calculated by $d_{\text{diff}} = d - R(a_1) - R(a_2)$, where $R(a_1)$ and $R(a_2)$ are the van der Waals radius of atom a_1 and a_2.^[47] q is the atomic charge and V is the atomic volume. ASP and QASP refer to atomic solvation parameter and charge-based solvation parameter, respectively.^[55] [b] The pharmacophore type definitions are presented in Supporting Information Table S3.</p>	

targets among all 65 targets. In CASF-2007, only the high-level success rate is evaluated.^[12]

Docking Power. The docking power test evaluates the ability of a scoring function to identify native binding poses among computer generated decoys. In CASF-2007, 100 decoy binding poses are selected from the poses generated by LigandFit,^[60] GOLD,^[61,62] Surflex,^[63,64] and FlexX^[65] for each ligand. Success is defined as one pose from the top one, the top two, or the top three poses ranked by predicted scores is within 2 Å RMSD from the native pose. The RMSDs of decoys relative to the native are provided in CASF-2007 benchmark and used directly. In CASF-2013, up to 100 decoy binding poses are selected from the poses generated by GOLD, Surflex and MOE. The property-matched RMSDs (RMSD^{PM}) of decoys, are calculated by considering the symmetry of the molecule in CASF-2013, however, RMSD^{PM} is not provided in CASF-2013. The symmetry-corrected RMSDs used here are calculated by Pybel.^[59] The native poses are included in the decoy set for success rate calculation.

Screening Power. The screening power test assesses the ability of a scoring function to identify a true binder from a pool of

random molecules for a given target. The test set for screening in CASF-2013 is designed by cross docking 195 ligands on 65 target proteins. For each protein, there are at least three true binders, as defined in Ref. 13. The remaining 192 ligands are searched through the ChEMBL database for possible cross-binders and 12 target proteins have more than 3 true binders in the dataset from the search. There are 12,675 (65 × 195) protein–ligand pairs from docking 195 ligands to 65 target proteins and up to 50 poses are selected for each protein–ligand pair. For a given target protein, 195 ligands are ranked based on the best-scored pose for a given protein–ligand pair. Screening power is measured by two metrics—enhancement factors and success rates—both based on the counts of the total number of true binders among the 1%, 5%, and 10% top-ranked molecules. Enhancement factors are computed for each target by

$$EF_{1\%} = \frac{NTB_{1\%}}{NTB_{\text{total}} \times 1\%}, EF_{5\%} = \frac{NTB_{5\%}}{NTB_{\text{total}} \times 5\%}, EF_{10\%} = \frac{NTB_{10\%}}{NTB_{\text{total}} \times 10\%}$$

$NTB_{1\%}$, $NTB_{5\%}$, and $NTB_{10\%}$ are the number of true binders among the 1%, 5%, and 10% top-ranked molecules. NTB_{total} is the total number of true binders for a given target protein.

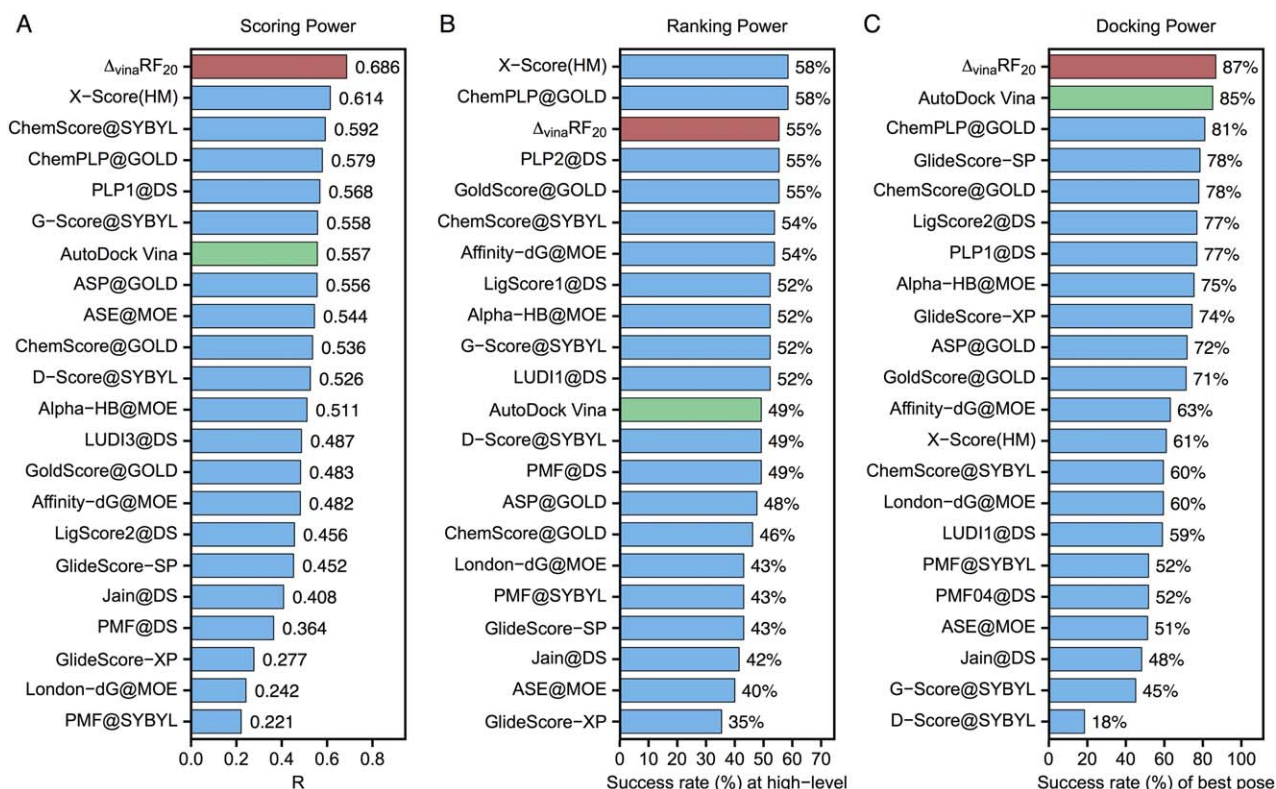


Figure 1. Performance of 22 scoring functions in (A) scoring power measured by Pearson's R, (B) ranking power in terms of high-level success rate, and (C) docking power measured by the success rate when the best-scored pose is considered to match the native pose in CASF-2013 benchmark. $\Delta_{\text{vina}}\text{RF}_{20}$ is colored in red and AutoDock Vina is colored in green. All results colored in blue are obtained from Ref. [13]. [Color figure can be viewed at wileyonlinelibrary.com]

The average EFs over 65 targets are calculated for each scoring function. Success rates are calculated as the number of targets that have true binders in 1%, 5%, and 10% of the top-ranked molecules among the total 65 targets.

Results

Based on a learning set which consists of 3336 experimental crystal structures and 3322 computer generated decoy structures, we have developed a new scoring function $\Delta_{\text{vina}}\text{RF}_{20}$ by employing random forest to parameterize a correction term to the original AutoDock Vina score.^[47] The overall scoring function $\Delta_{\text{vina}}\text{RF}_{20}$ has the following form: $pK_d(\Delta_{\text{vina}}\text{RF}_{20}) = pK_d(\text{Vina}) + \Delta pK_d(\text{RF}_{20})$, where $pK_d(\text{Vina})$ is the pK_d value calculated by multiplying the original Vina score by a factor of -0.73349 . $\Delta pK_d(\text{RF}_{20})$ is the correction term parameterized with random forest using 20 features as listed in Table 1. The $\Delta_{\text{vina}}\text{RF}_{20}$ model and code are available at: <https://www.nyu.edu/projects/yzhang/DeltaVina>. We have carried out all power tests on both CASF-2013 and CASF-2007 benchmark sets^[12,13] for $\Delta_{\text{vina}}\text{RF}_{20}$ as well as the AutoDock Vina scoring function,^[47] and compared with, respectively, 20 other and 16 other docking functions that were tested in the original 2013 and 2007 comparative assessment articles. The results are presented in Figures 1–3 and Supporting Information Tables S6–S13. The new scoring function $\Delta_{\text{vina}}\text{RF}_{20}$, which employs twenty descriptors in addition to the AutoDock Vina score, has achieved superior performance compared to classical scoring functions in

all tests of both CASF-2013 and CASF-2007 benchmarks, including scoring, ranking, docking, and screening power tests.

Scoring power

The scoring power of $\Delta_{\text{vina}}\text{RF}_{20}$ significantly outperforms AutoDock Vina as well as all scoring functions that have been tested in the original 2013 and 2007 comparative assessment articles, as shown in Figures 1 and 2. It achieves the best Pearson's correlation coefficients of 0.686 and 0.732 for the CASF-2013 and CASF-2007 benchmarks respectively, and significantly improves on AutoDock Vina, which has corresponding Pearson's correlation coefficients of 0.557 and 0.566, respectively.

Ranking power

The ranking power of $\Delta_{\text{vina}}\text{RF}_{20}$ is improved over AutoDock Vina, and is among the top 3 for both benchmarks. In CASF-2013, $\Delta_{\text{vina}}\text{RF}_{20}$ has a ranking power of 55% for high-level (in Fig. 1) and 74% for low-level (in Supporting Information Table S10), which places it third for high-level successes and second for low-level successes. In CASF-2007, the ranking power of $\Delta_{\text{vina}}\text{RF}_{20}$ is 57%, following the best X-Score::HSScore's success rate of 58%.^[66]

Docking power

The docking power of $\Delta_{\text{vina}}\text{RF}_{20}$ is among the top rank for both benchmarks. Its success rate to identify the top pose as

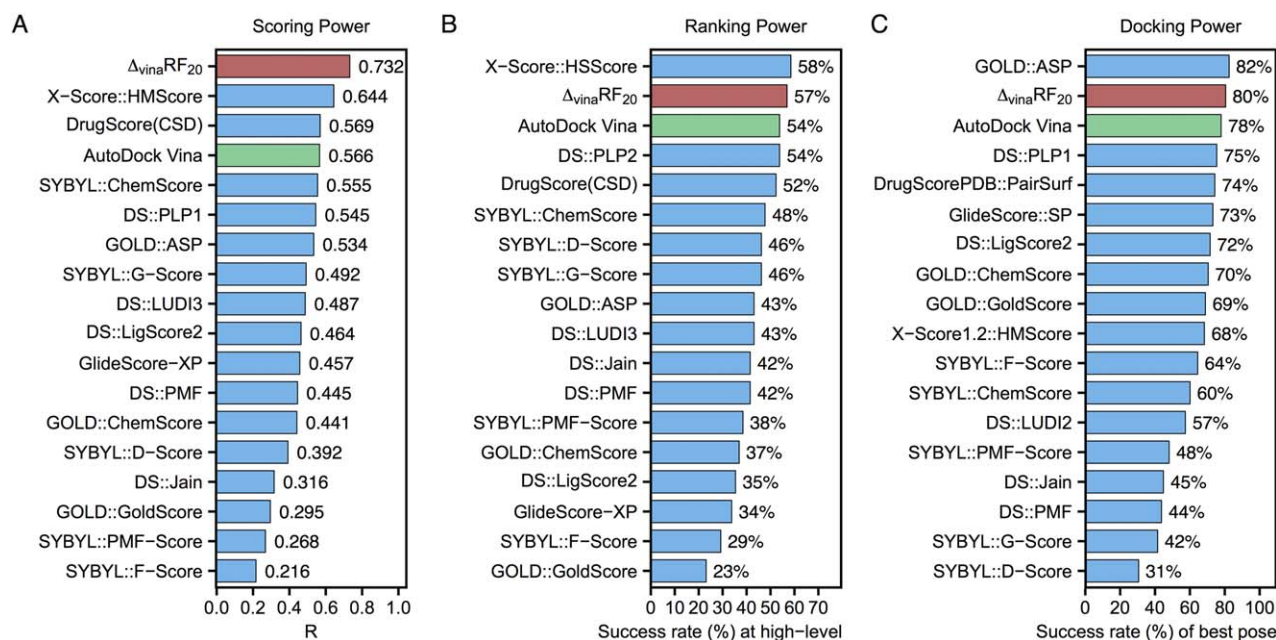


Figure 2. Performance of 18 scoring functions in (A) scoring power measured by Pearson's R, (B) ranking power in terms of high-level success rate, and (C) docking power measured by the success rate when the best-scored pose is considered to match the native pose in CASF-2007 benchmark. $\Delta_{\text{vina}}\text{RF}_{20}$ is colored in red and AutoDock Vina is colored in green. All results colored in blue are obtained from Ref. [12]. [Color figure can be viewed at wileyonlinelibrary.com]

the native pose is 87% in CASF-2013 (see Fig. 1), which outperforms all other scoring functions. In CASF-2007 (in Fig. 2), the success rate of $\Delta_{\text{vina}}\text{RF}_{20}$ is 80%, which ranks second following the best one GOLD::ASP (82%).^[67] We can see that the docking power of $\Delta_{\text{vina}}\text{RF}_{20}$ has improved on AutoDock Vina by about 2% in both benchmarks. Recently, a new SPA-SE score function was developed by combining knowledge-based atom-pair potential with the atomic solvation energy of a charge-independent implicit solvent model.^[68] It shows excellent performance in scoring (with a Pearson's correlation coefficient of 0.662), ranking (60.0% for high level and 75.4% for low level) and docking power (83.1% success rate for the best pose) for the CASF-2013 benchmark, while the screening power of SPA-SE was not reported. We can see that $\Delta_{\text{vina}}\text{RF}_{20}$ is still slightly better than SPA-SE in both scoring and docking power tests for the CASF-2013 benchmark.

Screening power

The screening power of $\Delta_{\text{vina}}\text{RF}_{20}$ is the best as shown in Figure 3 for both metrics: enrichment factor and success rate at top 1% level. The average enrichment factor of $\Delta_{\text{vina}}\text{RF}_{20}$ is 21 at top 1% level, which is slightly better than GlideScore-SP's 20, and the success rate is 60% (39 out of 65), which is the same as GlideScore-SP at top 1% level.^[19,20] It should be noted that $\Delta_{\text{vina}}\text{RF}_{20}$ has significantly improved on AutoDock Vina, which has an enrichment factor of 15.6 and success rate of 45% at top 1% level, respectively.

Discussion

Scoring functions play a central role in protein–ligand docking. An ideal, robust scoring function should perform well across

different important tasks, including scoring, docking, and screening power tests. Extensive retrospective and comparative studies^[10,11,14–18] indicate that although some widely used scoring functions can do relatively well in docking and screening power tests, most of them are weaker in performance in the scoring power test. Furthermore, it is very challenging for a docking function to achieve superior performance on all three power tests simultaneously.^[41,42] For example, X-Score(HM)^[66] is the top performer in the scoring power test in the original CASF-2013 comparative study, with a Pearson's correlation coefficient of 0.614, but its performance in the docking power test only ranks in the middle among about 20 tested scoring functions. Its success rate is 61% in predicting the best pose, which is significant lower than the value of 81% for ChemPLP@GOLD.^[69] These disparities in performance for different power tests become significantly worse for recently developed machine learning-based scoring functions (MLbScores).^[41,42] For example, in a recent comparative assessment of a dozen MLbScores^[42] for the CASF-2013 benchmark, RF@ML, which is parameterized with random forest using more than one hundred features, achieved the best scoring power of 0.704 in Pearson's correlation coefficient among all 12 scoring functions developed using different machine learning algorithms. However, the docking and screening powers of these 12 MLbScores are all significantly worse. The docking power of RF@ML is only 12.2% for success in predicting the best pose, while the screening power of RF@ML is just 6.45% for the success rate in finding the best ligand molecule. These values are significantly lower than those of classical scoring functions, whose top performances are 81% (ChemPLP@GOLD)^[69] and 60% (GlideScore-SP)^[19,20] for docking and screening power respectively. Conversely, as presented in the above results section, our newly developed $\Delta_{\text{vina}}\text{RF}_{20}$ scoring

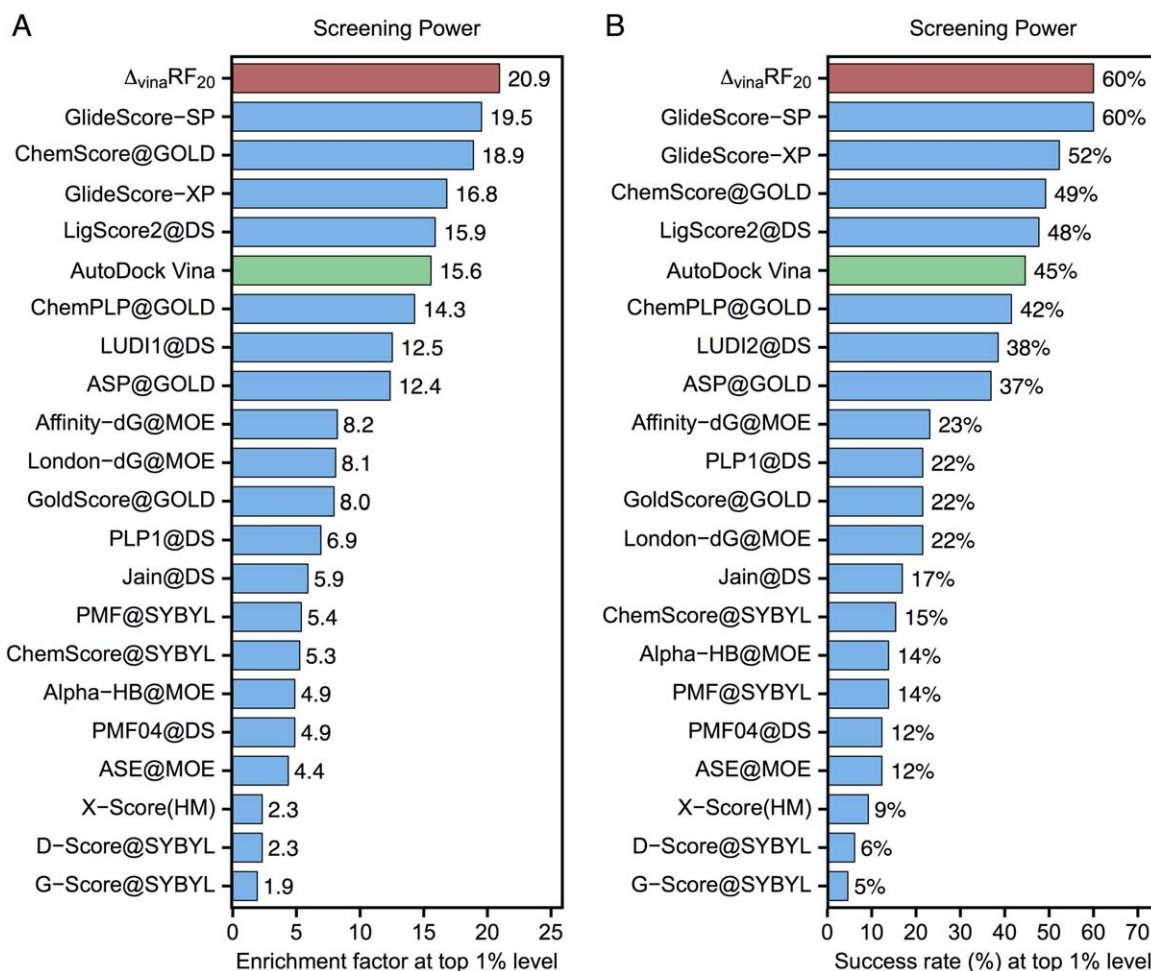


Figure 3. Performance of 22 scoring functions in screening power measured by (A) enrichment factor and (B) success rate at top 1% level in CASF-2013 benchmark. $\Delta_{vina}RF_{20}$ is colored in red and AutoDock Vina is colored in green. All results colored in blue are obtained from Ref. [13]. [Color figure can be viewed at wileyonlinelibrary.com]

function, using 20 features, achieved superior performance compared to traditional scoring functions in all power tests for both CASF-2013 and CASF-2007 benchmarks.

The main idea of our $\Delta_{vina}RF$ approach is to use random forest to parameterize a correction term to the AutoDock Vina score so that it can combine Vina's excellent docking power with RF's ability to significantly improve scoring accuracy. We tested both RF and $\Delta_{vina}RF$ approaches using the same twenty features and training set for the development of the $\Delta_{vina}RF_{20}$ scoring function, and compared them with Vina in Figure 4 for the scoring, docking, and screening power tests in CASF-2013. In addition, we have also trained a RF model using experimental data alone and the same twenty features. The results clearly demonstrate that the RF approach can only improve the scoring power by significantly sacrificing docking and screening powers, while the $\Delta_{vina}RF$ approach, with a combined experimental and decoy training set, can achieve the improvement over AutoDock Vina in all three tests simultaneously.

One attractive capability of the random forest algorithm is that it can efficiently utilize a large set of training data. It has been previously demonstrated that the larger the training data, the better the resulting RFbScore's performance in the

scoring power test,^[28,29,70] even when low-quality structures are included.^[30] For the $\Delta_{vina}RF$ approach, as shown in Supporting Information Figure S1, we also find that a larger experimental learning set can significantly improve the scoring power, but it does not necessarily improve the docking power. By expanding the learning set to include decoy structures, both docking and screening power of the scoring function can be improved over AutoDock Vina with the $\Delta_{vina}RF$ approach.

Besides the training set, another critical component of a random forest based scoring function (RFbScore) is the feature set. For a random forest model, both the number of features and feature relevance will affect its performance. Numerical experiments show that increasing the fraction of relevant features can improve the performance of the random forest model by increasing the chance that important features will be selected at each tree splitting.^[44] In this work, by taking advantage of random forest in ranking features, we employ a strategy that includes both feature selection and aggregation, to yield the 20 features in $\Delta_{vina}RF_{20}$. As listed in Table 1, the feature set of $\Delta_{vina}RF_{20}$ consists of five interaction terms and five ligand-dependent terms, which are selected from 58 terms implemented in the AutoDock Vina source code and 10 terms

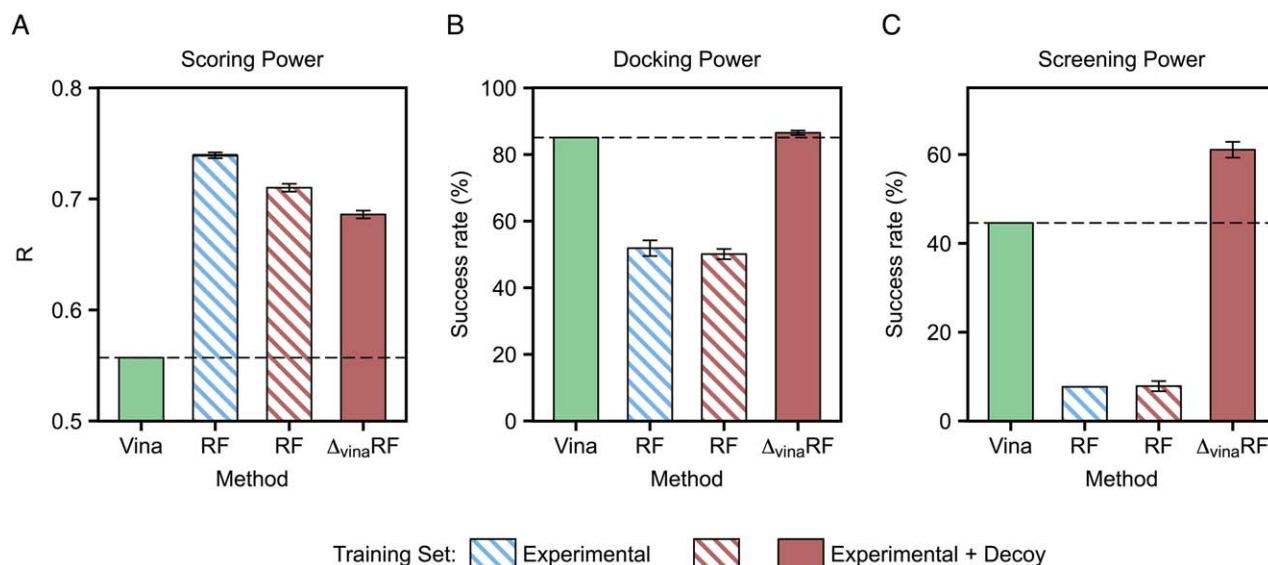


Figure 4. CASF-2013 benchmark test performance of AutoDock Vina (colored in green), scoring function developed with RF approach (colored in blue) using experimental data alone and the same twenty features in $\Delta_{vina}RF_{20}$ scoring function, and scoring functions developed with RF approach and $\Delta_{vina}RF$ approach (colored in red) using the same twenty features and the same training set for the development of the $\Delta_{vina}RF_{20}$ scoring function. (A) Scoring power; (B) Docking power; (C) Screening power. Each set is run 10 times with different random seed for random forest and calculated by averaging over 10 performances except AutoDock Vina. The AutoDock Vina performance is also indicated by dashed line. [Color figure can be viewed at wileyonlinelibrary.com]

related to buried solvent-accessible surface area (bSASA). Interestingly, all five interaction terms in the $\Delta_{vina}RF_{20}$ are related to polar interactions, and there is only one overlap term between $\Delta_{vina}RF_{20}$ and the original Vina score, which is the number of torsions in the ligand. In comparison with using only 6 terms in the original Vina score and all 58 terms implemented in the AutoDock Vina source code, the $\Delta_{vina}RF$ model, developed using 10 selected features, performs better in all scoring, docking and screening power tests for the CASF-2013 benchmark test as shown in Supporting Information Figure S2.

Among five ligand-dependent terms, rotors and torsions have been previously used to approximate the entropic change in several empirical scoring functions.^[19,66,71,72] Number of heavy atoms and ligand length could also be viewed as entropy-related features since both of them are highly correlated with rotors (Pearson's correlation coefficients are 0.802 and 0.906, respectively) in the crystal structure training set. Number of hydrophobic atoms is related to the hydrophobic interaction and the Pearson's correlation coefficient between number of hydrophobic atoms and the hydrophobic term defined in AutoDock Vina source code is around 0.85 for the crystal structure training set.

Surface area and related features are widely used in protein-ligand scoring function development due to their relation to solvation.^[66,68,71–74] The buried solvent-accessible surface area of a ligand (bSASA) has also been tested as a naive scoring function in CASF-2013 scoring comparison list^[13] and outperforms most of other scoring functions in the scoring power test but ranks the worst in both docking and screening power tests. Since no bSASA term has been implemented in the AutoDock Vina source code, we have explored nine pharmacophore-based bSASA terms and 1 total bSASA term as the feature set. As shown in Supporting Information Figure S3, the $\Delta_{vina}RF$ model

developed using 10 bSASA terms alone would also perform quite well in all three power tests, but not as good as when combined with 10 selected AutoDock Vina terms. By combining two feature sets, the resulting Vina10-bSASA, the feature set used in $\Delta_{vina}RF_{20}$, performs better than either Vina10 or bSASA in all three power tests for the CASF-2013 benchmark. We have tested the importance of features measured by percentage of increased mean squared error (%IncMSE) for the 20 features in $\Delta_{vina}RF_{20}$ as shown in Supporting Information Figure S4. Except for the halogen bSASA, which may be limited by the rarity of halogens in the training set, each of the other 19 features has a %IncMSE value significantly larger than 20%, indicating their general importance and justifying their inclusion in the feature set of $\Delta_{vina}RF_{20}$. Meanwhile, the additional test results in Supporting Information Figures S1–S4 further indicate the robustness of the $\Delta_{vina}RF$ approach.

Conclusion

A major challenge in developing a robust protein–ligand scoring function is to improve scoring, docking and screening performances simultaneously. In this work, we have made advances in overcoming this challenge by introducing a new $\Delta_{vina}RF$ parameterization and feature selection framework based on random forest. Our new scoring function $\Delta_{vina}RF_{20}$ employs twenty features in addition to the AutoDock Vina score, and can achieve superior performance compared to classical scoring functions in all tests of both CASF-2013 and CASF-2007 benchmarks, including scoring, ranking, docking and screening power tests. This work suggests that Δ -machine learning is a promising approach to systemically improve the performance and robustness of docking functions by

employing larger diverse experimental/decoy datasets of high quality, developing and selecting physically meaningful features, as well as adapting advanced machine learning algorithms.

Acknowledgment

We thank NYU-ITS and NYUAD for providing computational resources.

Keywords: random forest • docking • scoring function • protein–ligand binding affinity • machine learning

How to cite this article: C. Wang, Y. Zhang. *J. Comput. Chem.* **2017**, *38*, 169–177. DOI: 10.1002/jcc.24667



Additional Supporting Information may be found in the online version of this article.

- [1] S. Y. Huang, S. Z. Grinter, X. Q. Zou, *Phys. Chem. Chem. Phys.* **2010**, *12*, 12899.
- [2] P. D. Lyne, *Drug Discov. Today* **2002**, *7*, 1047.
- [3] B. K. Shoichet, *Nature* **2004**, *432*, 862.
- [4] C. McInnes, *Curr. Opin. Chem. Biol.* **2007**, *11*, 494.
- [5] R. V. C. Guido, G. Oliva, A. D. Andricopulo, *Curr. Med. Chem.* **2008**, *15*, 37.
- [6] T. J. Cheng, Q. L. Li, Z. G. Zhou, Y. L. Wang, S. H. Bryant, *AAPS J.* **2012**, *14*, 133.
- [7] A. Lavecchia, C. Di Giovanni, *Curr. Med. Chem.* **2013**, *20*, 2839.
- [8] D. L. Ma, D. S. H. Chan, C. H. Leung, *Chem. Soc. Rev.* **2013**, *42*, 2130.
- [9] R. X. Wang, Y. P. Lu, S. M. Wang, *J. Med. Chem.* **2003**, *46*, 2287.
- [10] P. Ferrara, H. Gohlke, D. J. Price, G. Klebe, C. L. Brooks, *J. Med. Chem.* **2004**, *47*, 3032.
- [11] P. M. Marsden, D. Puvanendrapillai, J. B. O. Mitchell, R. C. Glen, *Org. Biomol. Chem.* **2004**, *2*, 3267.
- [12] T. J. Cheng, X. Li, Y. Li, Z. H. Liu, R. X. Wang, *J. Chem. Inf. Model.* **2009**, *49*, 1079.
- [13] Y. Li, L. Han, Z. H. Liu, R. X. Wang, *J. Chem. Inf. Model.* **2014**, *54*, 1717.
- [14] I. Halperin, B. Y. Ma, H. Wolfson, R. Nussinov, *Proteins* **2002**, *47*, 409.
- [15] E. Perola, W. P. Walters, P. S. Charifson, *Proteins* **2004**, *56*, 235.
- [16] G. L. Warren, C. W. Andrews, A. M. Capelli, B. Clarke, J. LaLonde, M. H. Lambert, M. Lindvall, N. Nevins, S. F. Semus, S. Senger, G. Tedesco, I. D. Wall, J. M. Woolven, C. E. Peishoff, M. S. Head, *J. Med. Chem.* **2006**, *49*, 5912.
- [17] R. Kim, J. Skolnick, *J. Comput. Chem.* **2008**, *29*, 1316.
- [18] D. Plewczynski, M. Lazniewski, R. Augustyniak, K. Ginalska, *J. Comput. Chem.* **2011**, *32*, 742.
- [19] R. A. Friesner, J. L. Banks, R. B. Murphy, T. A. Halgren, J. J. Klicic, D. T. Mainz, M. P. Repasky, E. H. Knoll, M. Shelley, J. K. Perry, D. E. Shaw, P. Francis, P. S. Shenkin, *J. Med. Chem.* **2004**, *47*, 1739.
- [20] T. A. Halgren, R. B. Murphy, R. A. Friesner, H. S. Beard, L. L. Frye, W. T. Pollard, J. L. Banks, *J. Med. Chem.* **2004**, *47*, 1750.
- [21] R. A. Friesner, R. B. Murphy, M. P. Repasky, L. L. Frye, J. R. Greenwood, T. A. Halgren, P. C. Sanschagrin, D. T. Mainz, *J. Med. Chem.* **2006**, *49*, 6177.
- [22] M. I. Jordan, T. M. Mitchell, *Science* **2015**, *349*, 255.
- [23] Y. LeCun, Y. Bengio, G. Hinton, *Nature* **2015**, *521*, 436.
- [24] A. Lavecchia, *Drug Discov. Today* **2015**, *20*, 318.
- [25] M. W. Libbrecht, W. S. Noble, *Nat. Rev. Genet.* **2015**, *16*, 321.
- [26] P. J. Ballester, J. B. O. Mitchell, *Bioinformatics* **2010**, *26*, 1169.
- [27] P. J. Ballester, A. Schreyer, T. L. Blundell, *J. Chem. Inf. Model.* **2014**, *54*, 944.
- [28] H. J. Li, K. S. Leung, M. H. Wong, P. J. Ballester, *BMC Bioinf.* **2014**, *15*, 291.
- [29] H. J. Li, K. S. Leung, M. H. Wong, P. J. Ballester, *Mol. Inf.* **2015**, *34*, 115.
- [30] H. J. Li, K. S. Leung, M. H. Wong, P. J. Ballester, *Molecules* **2015**, *20*, 10947.
- [31] D. Zilian, C. A. Sotriffer, *J. Chem. Inf. Model.* **2013**, *53*, 1923.
- [32] Q. Liu, C. K. Kwok, J. Y. Li, *J. Chem. Inf. Model.* **2013**, *53*, 3076.
- [33] H. M. Ashtawy, N. R. Mahapatra, *BMC Bioinf.* **2015**, *16*(Suppl 4), S8.
- [34] J. D. Durrant, J. A. McCammon, *J. Chem. Inf. Model.* **2010**, *50*, 1865.
- [35] J. D. Durrant, J. A. McCammon, *J. Chem. Inf. Model.* **2011**, *51*, 2897.
- [36] I. Wallach, M. Dzamba, A. Heifets, *ePrint arXiv*, **2015**, arXiv:1510.02855. Available at: <http://arXiv.org/abs/1510.02855> (accessed on April 20, 2016).
- [37] L. Li, B. Wang, S. O. Meroueh, *J. Chem. Inf. Model.* **2011**, *51*, 2132.
- [38] B. Ding, J. Wang, N. Li, W. Wang, *J. Chem. Inf. Model.* **2013**, *53*, 114.
- [39] G. B. Li, L. L. Yang, W. J. Wang, L. L. Li, S. Y. Yang, *J. Chem. Inf. Model.* **2013**, *53*, 592.
- [40] W. Wang, W. L. He, X. Zhou, X. Chen, *Proteins* **2013**, *81*, 1386.
- [41] J. Gabel, J. Desaphy, D. Rognan, *J. Chem. Inf. Model.* **2014**, *54*, 2807.
- [42] M. A. Khamis, W. Gomaa, *Eng. Appl. Artif. Intel.* **2015**, *45*, 136.
- [43] L. Breiman, *Mach. Learn.* **2001**, *45*, 5.
- [44] T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*; Springer-Verlag, New York, **2009**.
- [45] A. J. Wyner, M. Olson, J. Bleich, D. Mease, Available at: <http://arXiv.org/abs/1504.07676> (accessed on April 20, 2016).
- [46] R. Ramakrishnan, P. O. Dral, M. Rupp, O. A. von Lilienfeld, *J. Chem. Theory Comput.* **2015**, *11*, 2087.
- [47] O. Trott, A. J. Olson, *J. Comput. Chem.* **2010**, *31*, 455.
- [48] Y. Li, Z. H. Liu, J. Li, L. Han, J. Liu, Z. X. Zhao, R. X. Wang, *J. Chem. Inf. Model.* **2014**, *54*, 1700.
- [49] R. X. Wang, X. L. Fang, Y. P. Lu, S. M. Wang, *J. Med. Chem.* **2004**, *47*, 2977.
- [50] R. X. Wang, X. L. Fang, Y. P. Lu, C. Y. Yang, S. M. Wang, *J. Med. Chem.* **2005**, *48*, 4111.
- [51] J. B. Dunbar, R. D. Smith, C. Y. Yang, P. M. U. Ung, K. W. Lexa, N. A. Khazanov, J. A. Stuckey, S. M. Wang, H. A. Carlson, *J. Chem. Inf. Model.* **2011**, *51*, 2036.
- [52] S. Y. Huang, X. Q. Zou, *J. Chem. Inf. Model.* **2011**, *51*, 2107.
- [53] A. Liaw, M. Wiener, *R News* **2002**, *2*, 18.
- [54] D. R. Koes, M. P. Baumgartner, C. J. Camacho, *J. Chem. Inf. Model.* **2013**, *53*, 1893.
- [55] R. Huey, G. M. Morris, A. J. Olson, D. S. Goodsell, *J. Comput. Chem.* **2007**, *28*, 1145.
- [56] M. F. Sanner, A. J. Olson, J. C. Spehner, *Biopolymers* **1996**, *38*, 305.
- [57] L. L. Jiang, R. C. Rizzo, *J. Phys. Chem. B* **2015**, *119*, 1083.
- [58] M. Clark, R. D. Cramer, N. Vanopdenbosch, *J. Comput. Chem.* **1989**, *10*, 982.
- [59] N. M. O'Boyle, C. Morley, G. R. Hutchison, *Chem. Cent. J.* **2008**, *2*, 5.
- [60] C. M. Venkatachalam, X. Jiang, T. Oldfield, M. Waldman, *J. Mol. Graphics Model.* **2003**, *21*, 289.
- [61] G. Jones, P. Willett, R. C. Glen, *J. Mol. Biol.* **1995**, *245*, 43.
- [62] G. Jones, P. Willett, R. C. Glen, A. R. Leach, R. Taylor, *J. Mol. Biol.* **1997**, *267*, 727.
- [63] A. N. Jain, *J. Med. Chem.* **2003**, *46*, 499.
- [64] A. N. Jain, *J. Comput. Aided Mol. Des.* **2007**, *21*, 281.
- [65] M. Rarey, B. Kramer, T. Lengauer, G. Klebe, *J. Mol. Biol.* **1996**, *261*, 470.
- [66] R. X. Wang, L. H. Lai, S. M. Wang, *J. Comput. Aided Mol. Des.* **2002**, *16*, 11.
- [67] W. T. M. Mooij, M. L. Verdonk, *Proteins* **2005**, *61*, 272.
- [68] Z. Q. Yan, J. Wang, *Proteins* **2015**, *83*, 1632.
- [69] O. Korb, T. Stutzle, T. E. Exner, *J. Chem. Inf. Model.* **2009**, *49*, 84.
- [70] H. M. Ashtawy, N. R. Mahapatra, *IEEE/ACM Trans. Comput. Biol. Bioinf.* **2015**, *12*, 335.
- [71] H. J. Bohm, *J. Comput. Aided Mol. Des.* **1994**, *8*, 243.
- [72] Y. Cao, L. Li, *Bioinformatics* **2014**, *30*, 1674.
- [73] A. Krammer, P. D. Kirchhoff, X. Jiang, C. M. Venkatachalam, M. Waldman, *J. Mol. Graph. Model.* **2005**, *23*, 395.
- [74] H. J. Bohm, *J. Comput. Aided Mol. Des.* **1998**, *12*, 309.

Received: 4 July 2016

Revised: 6 September 2016

Accepted: 26 October 2016

Published online on 17 November 2016