

Predicting Multiple Ligand Binding Modes Using Self-Consistent Pharmacophore Hypotheses

Izhar Wallach^{*,†,‡} and Ryan Lilien^{*,†,‡,§}

Department of Computer Science, Donnelly Centre for Cellular and Biomolecular Research, and Banting and Best Department of Medical Research, University of Toronto, Toronto, Ontario, Canada

Received June 3, 2009

The ability to predict ligand binding modes without the aid of wet-lab experiments may accelerate and reduce the cost of drug discovery research. Despite significant recent progress, virtual screening has not yet eliminated the need for wet-lab experiments. For example, after a lead compound has been identified, the precise binding mode is still typically determined by experimental structural biology. This structural knowledge is then employed to guide lead optimization. We present a step toward improving protein–ligand binding mode prediction for a set of ligands known to interact with a common protein. There is thus an important distinction between this work and traditional virtual screening algorithms. Whereas traditional approaches attempt to identify binding ligands from a large database of available compounds, our approach aims to more accurately predict the binding mode for a set of ligands which are already known to bind the target protein. The approach is based on the hypothesis that each active site contains a set of interaction points which binding ligands tend to exploit. In a more traditional context, these interaction points make up a pharmacophoric map. Our algorithm first performs traditional protein–ligand docking for each known binder. The ranked lists of candidate binding modes are then evaluated to identify a set of poses maximally self-consistent with respect to a pharmacophoric map generated from the same poses. We have extensively demonstrated the application of the algorithm to four protein systems (thrombin, cyclin-dependent kinase 2, dihydrofolate reductase, and HIV-1 protease) and attained predictions with an average RMSD < 2.5 Å for all tested systems. This represents a typical improvement of 0.5 – 1.0 Å (up to 25%) RMSD over the naive virtual docking predictions. Our algorithm is independent of the docking method and may significantly improve binding mode prediction of virtual docking experiments.

INTRODUCTION

Structural analysis of protein–ligand binding is fundamental to many components of the pharmaceutical development pipeline, including pharmacophore inference, high throughput screening (HTS), structure–activity relationship (SAR), and lead optimization.^{1–4} Native ligand binding modes are typically determined by solving the experimental structure of the protein–ligand complex using X-ray crystallography or nuclear magnetic resonance spectroscopy (NMR). These wet-lab techniques, however, are laborious, slow, and expensive. Docking algorithms aim to address these shortcomings by providing a fast and inexpensive technique for predicting binding modes. These computational methods have enjoyed varying degrees of success.^{5–7} Although docking algorithms are often useful for enriching small-molecule screening libraries,⁸ they can encounter significant challenge in predicting binding mode.^{5,9–11} In most cases, docking algorithms produce lists of candidate binding poses ranked by their predicted energetic favorability. Because energy functions only approximate the underlying biophysics, when near-native binding modes are identified by docking algorithms, they frequently fail to rank at the top of the lists.^{9,11}

We present a step toward improving protein–ligand binding mode prediction for a set of ligands known to interact with a target protein. We anticipate the method to be of use in HTS experiments where such families of binding ligands are known. Whereas traditional computational approaches attempt to identify binding ligands from a large database of available compounds, our approach aims to more accurately predict the binding mode for a set of ligands already known to bind the target protein. For example, HTS experiments may identify several ligands capable of binding a target protein. These ligands are expected to have similar (though not necessarily identical) interaction patterns (i.e., a common underlying pharmacophoric map). The problem of identifying a pharmacophoric map using only active ligands without structural knowledge of the protein–ligand complex is known as ligand-based pharmacophore inference. Algorithms for ligand-based pharmacophore inference search for a map that is maximally consistent with the binding of all ligands to the protein.^{12–14} In this scenario, the accuracy of the pharmacophoric map is likely to correlate with the correctness of the binding modes. This suggests that the search can be reversed. It may be possible to identify native binding modes by selecting poses consistent with this pharmacophoric map. Our approach exploits this hypothesis. Starting with a set of ligands known to bind a target protein, we generate ranked lists of candidate binding modes using a traditional docking algorithm. These poses are then evaluated to identify

* Corresponding author e-mail: izharw@cs.toronto.edu (I.W.); lilien@cs.toronto.edu (R.L.).

† Department of Computer Science.

‡ Donnelly Centre for Cellular and Biomolecular Research.

§ Banting and Best Department of Medical Research.

a subset of poses maximally self-consistent with respect to an underlying pharmacophoric map. The pharmacophoric map is not known *a priori* but is iteratively generated and refined during the search. We propose that the accuracy of predicted binding modes can be evaluated by scoring the underlying pharmacophoric map. Such a score would reflect the agreement of the ligands over the locations and chemical labels of the pharmacophoric points inside the binding site. We devise an objective function such that a map generated from a set of near-native binding modes will score higher than a map generated from non-native binding modes. Under this assumption, identifying the highest scoring pharmacophoric map implies finding a set of near-native binding modes.

Several techniques have been developed to utilize ligand similarity for binding mode prediction. Pharmacophoric maps generated from sets of binding ligands have successfully aided virtual docking. These maps focus on chemical interaction patterns and topology observed within the protein binding site.^{15,16} Maps generated from sets of holo-complexes have been particularly successful at discriminating active ligands from decoys in virtual screening experiments for protein kinase inhibitors.¹⁷ Evaluating candidate ligands using these maps efficiently discriminates decoys and yields binding mode similar to native poses. However, this method requires an *a priori* pharmacophoric model of the target protein. Thus, it is limited to protein targets for which there is an available set of previously solved holo complexes. Another approach for binding mode prediction utilizes a user defined scaffold common to all ligands.¹⁸ The scaffold serves as an anchor reference within the protein binding cavity. Possible binding modes are selected or eliminated according to their fit on the scaffold.¹⁹ This approach yields superior binding mode predictions than docking each ligand individually but is limited to sets of highly similar ligands sharing a common scaffold. Renner et al. recently introduced a binding mode prediction method using a Maximum Common Binding Modes (MCBM) approach.²⁰ Their method uses structural interaction fingerprints (SIFT) to model protein–ligand interactions.²¹ An exhaustive search over triplets of interaction fingerprints identifies the mode that maximally satisfies a set of interactions common to all ligands. While this method is promising, it may be sensitive to the degree of similarity within the ligand set – prediction quality degraded when the ligand similarities decreased. In addition, the method’s focus on the interface between the protein and the binding ligands might ignore structural information such as scaffold similarity. This structural knowledge is not only useful for binding mode prediction¹⁸ but also provides geometrical constraints essential for a reliable inference of protein–ligand interaction patterns.²² Finally, binding patterns of ligands not directly interacting with the protein (i.e., interaction via a water molecule) may not be detected by the SIFT method, even though a corresponding pharmacophoric point will be present.

We suggest a new hybrid approach for both predicting binding modes and simultaneously generating a pharmacophoric map. First, we use virtual docking to generate and rank possible binding modes. Then our algorithm identifies a set of poses maximally consistent with their inferred underlying pharmacophoric models. In contrast to previous pharmacophore-based approaches,^{16,17} our pharmacophoric

maps do not need to be known *a priori*. Our meta-analysis approach resembles consensus scoring techniques²³ where different scoring systems are applied over the same input. However, in our case the algorithm uses a single scoring system to explore multiple independently docked ligands while employing meta-data, such as chemical group similarities and spatial overlap, for the binding mode inference process.

In the following, we establish the ability to discriminate between candidate binding modes using our pharmacophore-based approach. We demonstrate the strength of our algorithm by accurately predicting binding modes for four protein systems over a large range of prediction experiments. The robustness of the algorithm is evaluated with respect to the number of ligands in the docking set, the similarity of the ligands, the number of correct binding modes generated by the docking algorithm, and the number of decoy ligands in the set (the term decoy is defined in “Prediction with Decoys” below).

METHODS

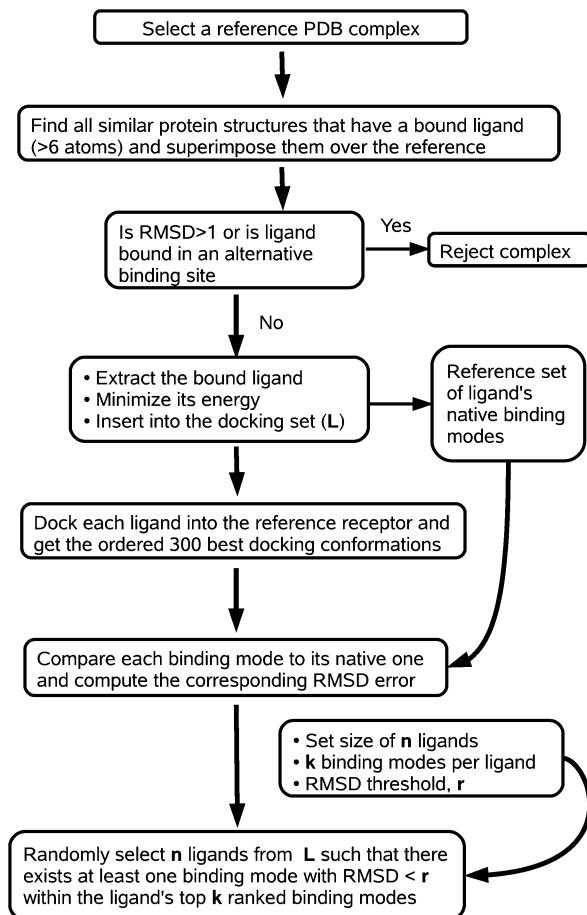
We present an algorithm to identify a maximally consistent set of binding modes for a set of ligands known to bind a target protein. Our method utilizes the spatial arrangement of the ligands’ chemical groups to first generate a pharmacophoric map and then discriminate between candidate binding modes. We note that any virtual docking algorithm capable of generating ranked lists of binding modes can be used with our technique. This requirement allows our method to be used within a range of screening environments.

In this paper, we have divided our methods into two parts. The first part (Figure 1a) generates input sets of virtual docking results, whereas the second part (Figure 1b) applies the binding mode prediction algorithm over these sets. The binding mode prediction algorithm is the focus of this work, and we expect it to be of interest to the drug discovery community. For completeness, we describe the methods used to generate the data sets employed in testing; we stress the importance of data set generation to facilitate an unbiased and thorough assessment of the prediction results.

Data Set Generation. An ideal molecular modeling algorithm will perform consistently well across a range of different protein systems. To demonstrate the consistency of our method, we generate a large testing set and perform a range of experiments with varying assumptions and quality of the candidate binding modes.

Given a holo protein target we define its binding site as the set of all amino-acid residues within 15 Å of a ligand atom. We search the Protein Data Bank²⁴ for other holo-structures which meet the following criteria:

1. The bound ligand has more than six non-hydrogen atoms.
 2. The ligand binds the protein at the binding site defined by the reference complex (i.e., within 5 Å of a reference ligand’s atom).
 3. The protein structure can be aligned to the target protein with an RMSD < 1 Å for non-hydrogen atoms.
- These rules identify sets of ligands capable of binding the target protein as well as the correct binding mode for each compound. We align the selected complex to the target protein, extract the ligand from the complex, and keep a copy



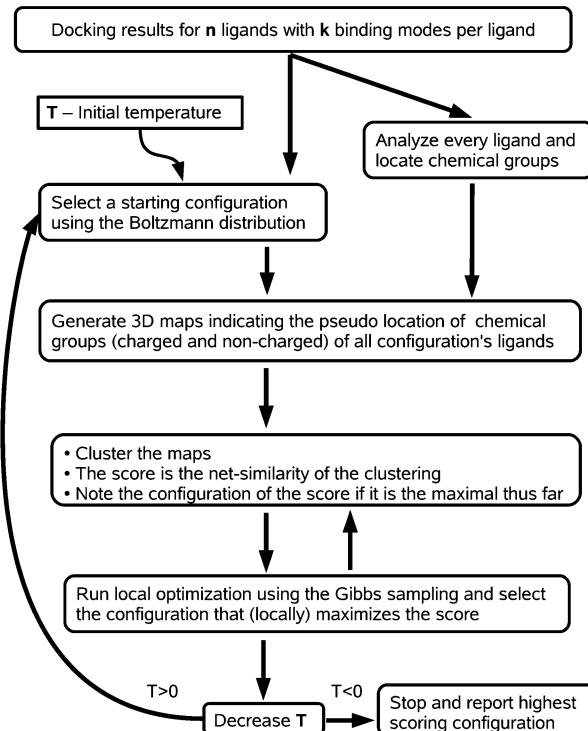
(a) Dataset preparation

Figure 1. (a) Data set generation - These methods identify sets of ligands known to bind a target protein and then generate (via docking) lists of candidate binding modes for each compound. These compounds served as the test cases described in this manuscript. The experimentally determined binding modes of each PDB structure were used as references when computing RMSDs. (b) Binding mode prediction. These methods identify sets of binding modes (configurations) that are maximally consistent with their underlying pharmacophoric model.

of the ligand as a ‘native binding mode’ reference (i.e., the *correct* binding mode for that ligand). The potential energy of every ligand is minimized, and the ligand is added to the input virtual docking set, L . We dock all ligands in L onto the target protein and obtain, for each ligand, a list of the top 300 binding modes ranked by the docking algorithm. For each ligand in L , we note the RMSD error of the ranked binding modes by comparing each docked pose to the native binding mode.

During evaluation of our algorithm we consider data sets of n ligands each associated with k candidate binding modes. We generate inputs for these experiments by sampling the docked poses of L above. For example, generating a data set of n ligands, each having k binding modes, is done by randomly selecting n ligands from L and using the first k ranked binding modes ($k \leq 300$) of each ligand. In order to generate sets where every ligand has at least one correct binding mode (i.e., a mode closer than 2.5 Å RMSD to the native) in its ranked list, we construct a set \tilde{L} of all ligands in L that have at least one binding mode for which the RMSD is smaller than a given threshold. Then, we use \tilde{L} instead of L to select lists of binding modes.

Binding Mode Prediction. The binding mode prediction algorithm receives an input set of ligands and their ranked lists of candidate binding modes. Our method searches for a



(b) Configuration prediction

set of binding modes (one per ligand) that is maximally self-consistent with its underlying pharmacophoric map. This is achieved by maximizing an objective function defined in such a way that large scores correspond to configurations with low RMSD to the (typically unknown) native binding mode.

For clarification, we define a *configuration* to be a set of poses, each pose belonging to a different ligand. Predicting the binding mode for a set of ligands corresponds to finding a configuration for which all poses are near-native. The RMSD of a pose is defined with respect to its native binding mode, whereas the RMSD of a configuration is the average RMSD of the configuration’s poses. Formally we define the binding mode prediction problem as the following:

Let L_n^k be a set of n ligands where each ligand has k poses indexed such that $\vec{l}_i^j \in L_n^k$ is the j -th pose of the i -th ligand and $1 \leq i \leq n$, $1 \leq j \leq k$. We define a configuration, C_n^k , to be a set of n binding modes, one from each ligand, such that $C_n^k = (l_1^{c_1}, \dots, l_n^{c_n})$ where $1 \leq c_i \leq k$ is the index of the pose of the i -th ligand. Given a scoring function f , which evaluates the consistency of a configuration, we look for a configuration such that

$$\arg \max_{C_n^k} f(C_n^k)$$

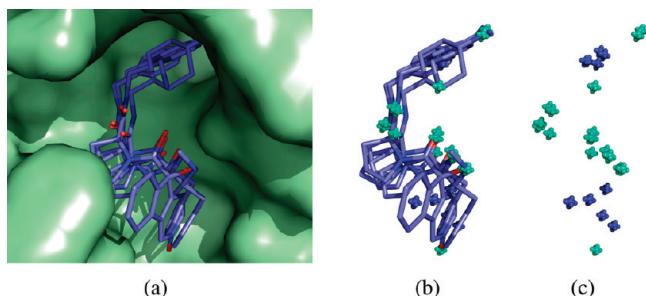


Figure 2. The generation of 3-dimensional functional maps: (a) superposition of 5 thrombin holo complexes (1TOM, 1ZRB, 2ZDA, 2ZFP, and 3BIV) using 1TOM as a reference, (b) the superimposed ligands with the location of identified functional groups (indicated with +), and (c) a 3-dimensional map of both polar (green) and nonpolar (blue) functional groups.

Below we elaborate the steps of the prediction algorithm as illustrated in Figure 1b.

Chemical Groups Analysis. In the context of this paper, we define the following terms. A *chemical group* is a group of atoms that characterizes a chemical moiety. We utilize a set of 47 chemical groups (e.g., phenyl, hydroxyl, carboxyl) inspired by Chen et al.²⁵ A *function type* (or a *functional group*) specifies the chemical functionality of a given group. We utilize six canonical function types: hydrophobic, aromatic, acid, base, hydrogen-bond donor (HBD), and hydrogen-bond acceptor (HBA).¹⁷ As some chemical groups correspond to more than one function type (i.e., HBD and HBA), the mapping between a chemical group and functional types can be one-to-many. In practice, each of our 47 chemical groups can be described by one of nine different sets of function types. We identify chemical groups over the ligands by using a template library encoded using SMARTS patterns²⁶ and then match these templates against each ligand structure. Similar to Schmitt et al.,²⁷ we assign a pseudolocation for each identified chemical group and label the group with up to six features (e.g., HBD, HBA,...) corresponding to its function type. We note that it is often difficult to differentiate acid from base and HBA from HBD without sophisticated modeling of the context of the receptor. Fortunately, as will be explained in the following section, our algorithm is invariant with respect to this possible ambiguity.

3D Maps of Functional Group Distribution. Given a configuration, C_n^k we generate two 3-dimensional maps containing the location and type of functional groups of each ligand $l_i^j \in C_n^k$ (Figure 2). One functional map indicates the locations of polar functional groups (acid, base, HBD, HBA), and the second map indicates the locations of the nonpolar groups (hydrophobic, aromatic). By using two types of functional maps rather than six, we attribute a greater importance to the spatial overlap of functional groups as compared to the actual match of their chemical functions. The polar/nonpolar representation allows the algorithm to match chemically close groups (e.g., HBD and acid) and to disallow matching where functional groups are not likely to be in agreement over a pharmacophoric point.

Clustering and Scoring Pharmacophoric Maps. A pharmacophoric-like map can be computed by identifying clustered functional groups within the two maps. Given the polar and nonpolar maps, the Affinity Propagation algorithm²⁸ is used to cluster and score the functional groups.

Affinity propagation requires the pairwise similarity to be specified between pairs of functional groups in each map. We define the similarity between two functional groups, p_i and p_j , as the negative distance between the groups using one of two distance functions:

Euclidean Similarity. $S(p_i, p_j) = -\|p_i - p_j\|_2$. This similarity function imposes a linear correlation between distance and similarity. That is, two pairs of functional groups, the second at twice the distance of the first will be scored half as similar.

Sigmoid Similarity. $S(p_i, p_j) = -1/(1 + \exp(c - \|p_i - p_j\|_2))$ where c is a parameter that defines the point at which the Sigmoid function reaches 50% similarity with respect to the Euclidean distance between the two groups. The Sigmoid similarity function defines a nonlinear correlation relating the Euclidean distance between a pair of groups and their similarity.

The Affinity Propagation algorithm defines a net similarity score, which is the sum of all similarities between each cluster member and the cluster's exemplar plus the sum of all exemplar preferences.²⁸ The net similarity score is an indicator of how well the objective function has been maximized and in our case reflects the quality of the clustering. Given the net similarities, we define the score of a configuration to be the sum of net similarities within both the polar and nonpolar maps. We hereafter refer to this score as the *configuration score*.

Optimization. An exhaustive configuration search among all combinations of binding poses is generally infeasible, even for a small set of ligands. For example, a set of 7 ligands, each having a ranked list of 100 binding modes, yields 10^{14} possible configurations. We address this combinatorial optimization problem by repeatedly applying local optimization steps over different initial configurations.

Choosing Initial Configurations. We direct the selection of the initial configurations toward binding modes for which we have prior information, indicating their likeliness to be correct. Using the virtual docking score, we assign each pose a probability that expresses our prior belief in its correctness. Then we sample initial configurations from these probability distributions. The probability of selecting a binding mode l_i^j is defined by the Boltzmann distribution

$$P(l_i^j) = \frac{\exp\left(-\frac{E_i^j}{k_B T}\right)}{\sum_{m=1}^k \exp\left(-\frac{E_i^m}{k_B T}\right)}$$

where E_i^m is the virtual docking score of the m -th binding mode (i.e., pose) of the i -th ligand, k_B is the Boltzmann constant, and T is the system temperature. Initially we use a high temperature, thereby allowing a nearly uniform sampling of configurations. Then, we gradually cool the system to move the solution toward configurations with high virtual docking scores. We sample 20 initial configurations at each temperature and use a cooling schedule of 0.95. Each initial configuration is optimized using Gibbs sampling.

Local Maximization Using Gibbs Sampling. Given a starting configuration, C_n^k we use Gibbs sampling²⁹ to search for a new configuration with an optimal score. We order the n ligands by their largest virtual docking score (ordered worst to best). Then, we iterate over the ligands in that order (i.e., worst to best) and compute the configuration score of the $k - 1$

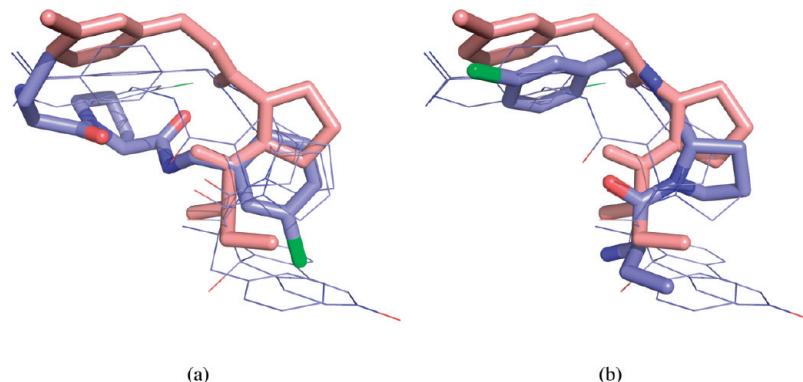


Figure 3. An example of one Gibbs optimization step for a set of thrombin ligands (2ZDA, 3BIV, 1ZRB, and 2ZFP). (a) Given an initial configuration, the algorithm evaluates all candidate binding modes for the current ligand (2ZFP, thick blue wireframe) with respect to the binding modes of the other ligands in the configuration (thin blue wireframe). (b) The highest scoring binding mode (thick blue wireframe) is selected for the current ligand (for reference, the native binding mode is shown in pink).

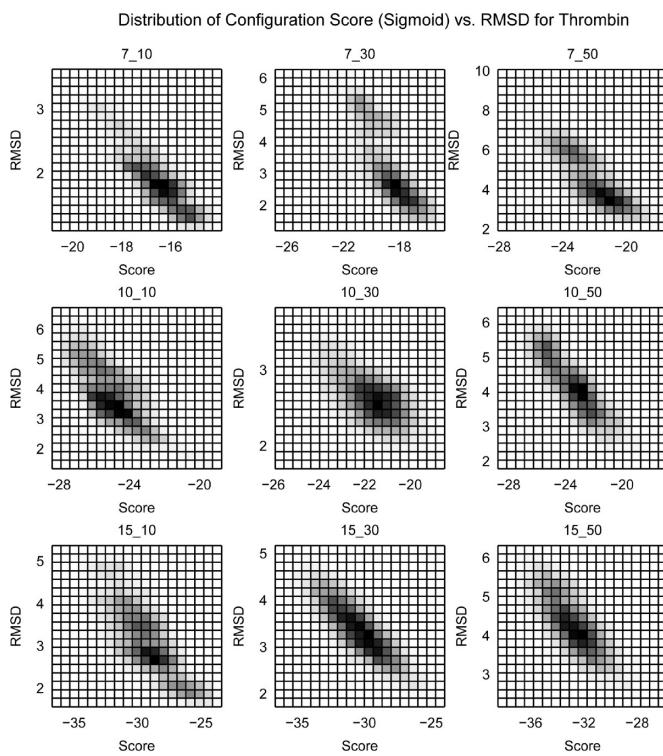
possible configurations obtained by substituting the pose of the current ligand with any of its other $k - 1$ possible poses. We select the configuration that maximizes the current configuration score and proceed to the next ligand (Figure 3). The algorithm terminates when a complete iteration over all ligands does not identify an improved configuration. We note that this optimization approach converges to a local maximum since we use steepest-ascent steps to move within the configuration space. The overall process terminates when the temperature of the system reaches zero.

RESULTS AND DISCUSSION

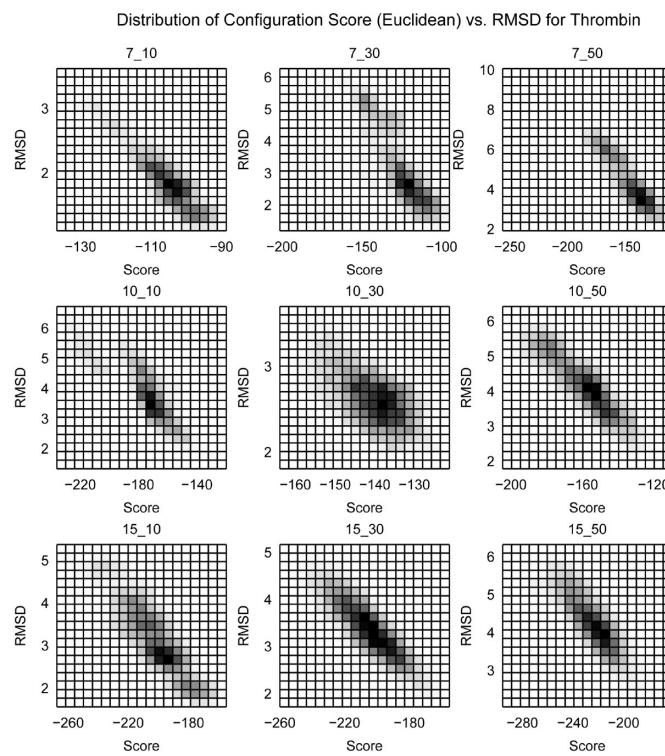
We performed two series of experiments. In the first, we established the ability of the scoring function (i.e., the

configuration score) to discriminate ligand binding modes. Then, we demonstrated the application of our algorithm using four protein target systems - thrombin, cyclin-dependent kinase 2 (CDK2), dihydrofolate reductase (DHFR), and HIV-1 protease (HIV-1P). The availability and diversity of holo complexes of these systems in the PDB allowed us to generate large data sets of ligands with known native binding modes.

We selected the following PDB complexes as protein models: 1TOM, 1OIT, 1BOZ, and 1AJV for thrombin, CDK2, DHFR, and HIV-1P, respectively. We obtained from the PDB 94, 114, 67, and 174 holo complexes for the thrombin, CDK2, DHFR, and HIV-1P docking sets (i.e., set L introduced in section “Data Set Generation”). Using the



(a) Sigmoid



(b) Euclidean

Figure 4. Two-dimensional histograms showing the distribution of the configuration score vs RMSD of 10^4 random configurations using the (a) Sigmoid ($c = 3$) and (b) Euclidean similarity functions and thrombin (1TOM) as a target protein. The title of each subplot, n_k , indicates a set of n ligands each having a ranked list of length k . No significant change in quality was observed for $c = \{3\ldots7\}$ (results not shown), and a value of 3 was selected for the remainder of this manuscript.

Table 1. Pearson Correlation Coefficient of Configuration Score vs RMSD Using the Sigmoid and Euclidean Similarity Functions and Thrombin (1TOM) as the Target Protein^a

set size (<i>n</i>)	no. of binding modes (<i>k</i>)	Pearson correlation	
		Sigmoid (<i>c</i> = 3)	Euclidean
7	10	-0.84	-0.89
7	30	-0.74	-0.81
7	50	-0.80	-0.89
10	10	-0.73	-0.83
10	30	-0.59	-0.60
10	50	-0.84	-0.88
15	10	-0.77	-0.85
15	30	-0.83	-0.86
15	50	-0.75	-0.82
mean		-0.77	-0.82

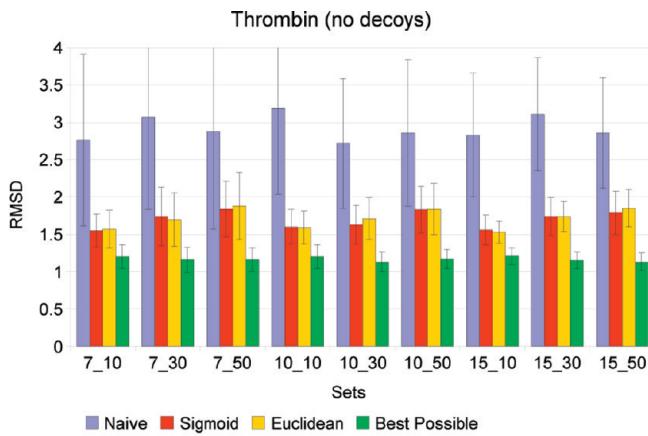
^a p-value of all correlations are <10⁻¹⁰.

FlexX docking package³⁰ (version 3.1.2), we generated ranked lists as described in the “Data Set Generation” section. The docking was performed with full ligand flexibility and a rigid protein under the default parameters. Omitting ligands that had no binding mode pose with RMSD < 2.5 Å left us with sets of 57, 49, 33, and 49 ligands (\tilde{L}).

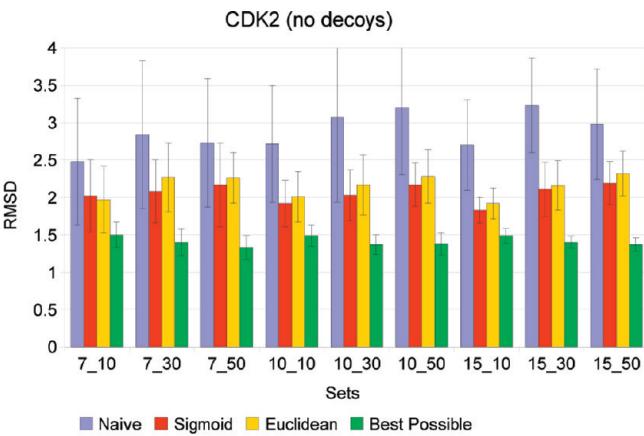
Scoring Function Analysis. We first established the ability of the configuration scoring function to find optimal configurations by showing that a configuration’s score correlates with its RMSD. Then, we evaluated the configuration scoring function using Sigmoid and Euclidean similarity measurements under varied sets of input ligands.

Correlation between Score and RMSD. In the general operating setting, the native binding modes are unknown and therefore cannot be used to score the quality of a configuration. Our configuration scoring function serves as a surrogate for the RMSD in scoring configurations. The scoring function was defined under the hypothesis that the poses of a correct configuration will be maximally consistent over an underlying pharmacophoric map. Consequently, the configuration score should correlate with the configuration’s RMSD. In the first set of experiments we test this hypothesis and find that configurations with a favorable score have favorable RMSDs.

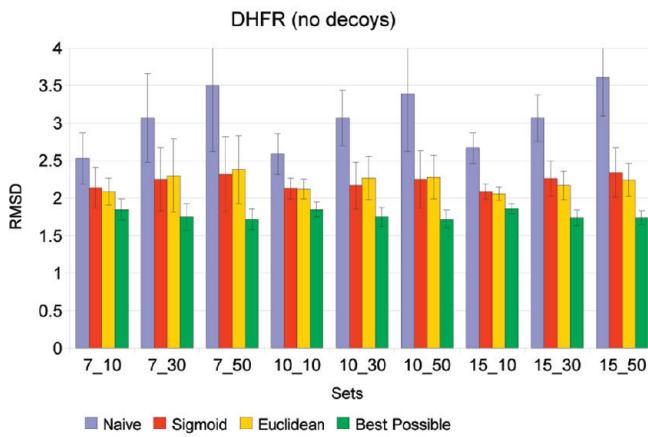
Three parameters can influence the relationship between a configuration’s score and its RMSD: (i) the number of ligands in the docking set, which affects the number of functional group clusters and their sizes, (ii) the similarity



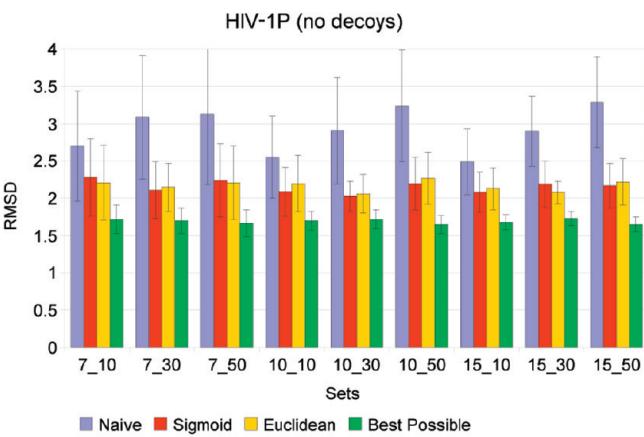
(a)



(b)



(c)



(d)

Figure 5. RMSD prediction results for systems without decoys for thrombin, CDK2, DHFR, and HIV-1P using varied set sizes and length of ranked lists. A set, n_k , implies a set of experiments where n ligands each have a ranked list of k possible binding modes. We compare predictions using the Sigmoid (red) and Euclidean (yellow) similarity functions to the naive approach (blue) and the best-possible configuration (green). Each bar indicates the average RMSD over 50 random experiments using the corresponding n_k parameter setting (“Data Set Generation” section).

of the ligands, which affects the distribution of chemical groups in space and thus the compactness of the clusters, and (iii) the functional group similarity function which affects the distribution of the clusters in space. We evaluated the configuration scoring function under different sets of these parameters. For each set of parameters, we randomly sampled 10^4 configurations and evaluated the resulting score and RMSD.

Figure 4 shows 2-dimensional histograms of score and RMSD for these random configurations using thrombin as the target protein. Similar results were obtained for different random sets of parameters and different targets (results not shown). The figure demonstrates a negative correlation between configuration score and RMSD with high tolerance to the tested parameters (i.e., configurations with a higher score have a lower RMSD). Table 1 shows the Pearson correlation coefficients for the score vs RMSD. We observe significantly strong negative correlations that remain high, regardless of the size of the sets, the length of the ranked lists, or the similarity function used. Furthermore, the correlations obtained using the Euclidean similarity function are moderately better than those obtained using the Sigmoid function. We note that the lower linear correlation obtained

using the Sigmoid function does not necessarily imply inferior predictive power. In fact, nonlinear correlation between score and RMSD might yield better results if the RMSD value decreases faster than the increase of the configuration's score. In other words, we do not necessarily require a linear correlation between score and RMSD, only that as the score increases the RMSD decreases.

Predictions. We demonstrated the application of our algorithm over ligand sets known to bind to thrombin, CDK2, DHFR, and HIV-1P. We first evaluated sets for which every ligand had at least one correct binding pose in its ranked list. Then, we evaluated sets with an increasing number of decoys, where we defined a decoy to be a ligand that did not contain a correct binding mode in its ranked list (RMSD $< 2.5 \text{ \AA}$).

Prediction When Every Ligand Has a Correct Binding Mode. We evaluated the ability of the algorithm to select optimal configurations in the ideal case where every input ligand has at least one correct binding pose in its ranked list. We ran 50 random prediction experiments (Figure 1b) for each parameter setting using the two similarity functions and varied sizes of the input sets and length of ranked lists (7, 10, and 15 ligands with 10, 30, and 50 poses). Figure 5 summarizes the prediction

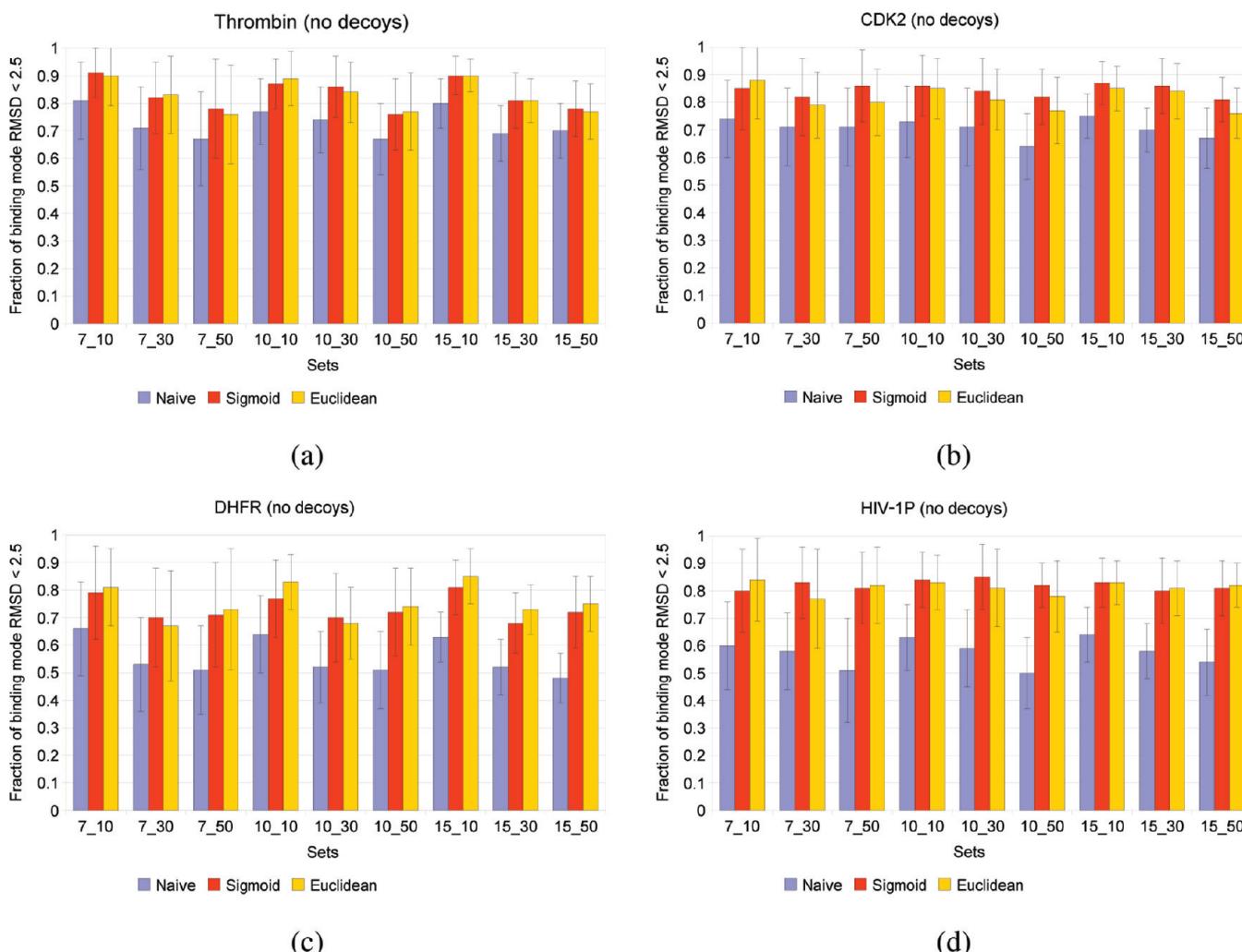


Figure 6. Fraction of correct binding modes for systems without decoys for thrombin, CDK2, DHFR, and HIV-1P using varied set sizes and length of ranked lists. A set, n_k , implies a set of experiments where n ligands each have a ranked list of k possible binding modes. We compare predictions using the Sigmoid (red) and Euclidean (yellow) similarity functions to the naive approach (blue). Each bar indicates the average fraction of correct binding modes (RMSD $< 2.5 \text{ \AA}$) over 50 random experiments using the corresponding n_k parameter setting (“Data Set Generation” section).

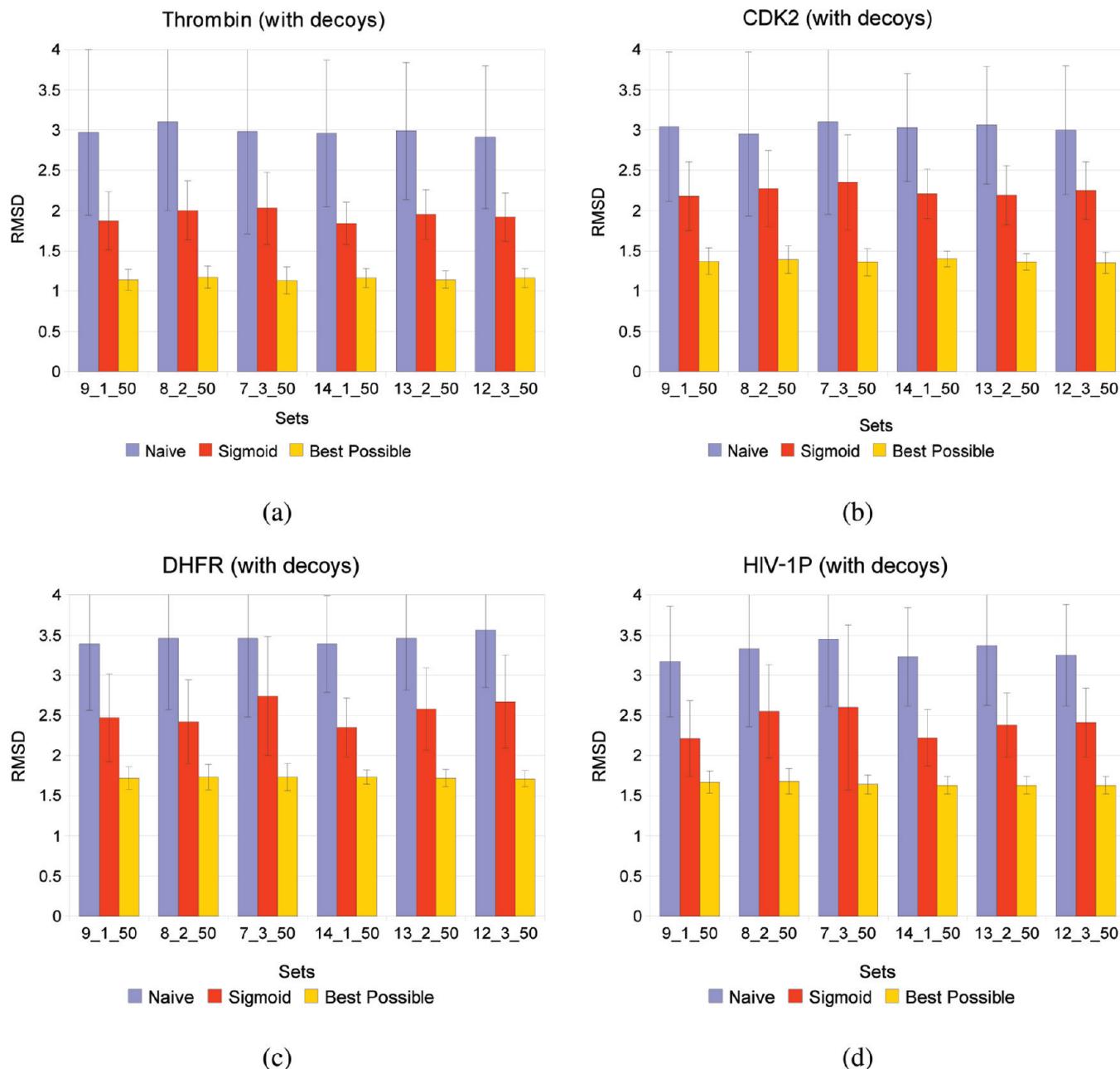


Figure 7. RMSD prediction results for systems with 1, 2, and 3 decoys for thrombin, CDK2, DHFR, and HIV-1P using varied set sizes and length of ranked lists. A set, x_y_k , implies a set of experiments where x ligands have correct binding modes, y ligands are decoys, and each ligand has a ranked list of k possible binding modes. Comparison between the average RMSD of the predicted sets (red) and the naive prediction (blue). The best possible configuration is shown as a reference (yellow). Each bar indicates the average RMSD over 50 random experiments using the corresponding x_y_k parameter setting (“Data Set Generation” section).

results for the four target system with respect to the *naive* and *best possible* configurations, using Euclidean and Sigmoid similarity functions (3600 total prediction experiments). The naive configuration consists of the highest ranked pose of every ligand in the set (i.e., the pose with the highest docking score). The best possible configuration consists of the ideal case where the pose closest to the native is selected for every ligand in the input set. The new algorithm shows significant improvement in the configuration selection over the naive approach (Figure 5). The average RMSD of the predicted binding modes was below 2.5 Å for all protein systems under all tested parameters. Notably, the thrombin results were even better, with an RMSD below 2 Å. The improvements over the naive approach were 1.22 Å, 0.83 Å, 0.84 Å, and 0.77 Å for the thrombin, CDK2, DHFR, and HIV-1P, respectively. The results also suggest that

the Euclidean and Sigmoid similarity functions are equally effective. We note that the shorter the ranked list is, the better the naive predictions become. This phenomenon results from the data set generation process which guarantees the presence of a correct binding mode for every input ligand in its k -ranked list. Ligands satisfying this constraint are likely to be easier to dock, and because the lists are shorter, are more likely to have a correct binding mode ranked first (the section “Robustness of the Algorithm” provides an analysis of the binding mode distribution). Figure 6 illustrates the fraction of poses in the predicted configuration that have $\text{RMSD} < 2.5 \text{ \AA}$. The algorithm demonstrates relative (and absolute) percent improvements of 14%(10%), 19%(14%), 32%(18%), and 43%(25%) over the naive predictions for the thrombin, CDK2, DHFR, and HIV-1P, respectively.

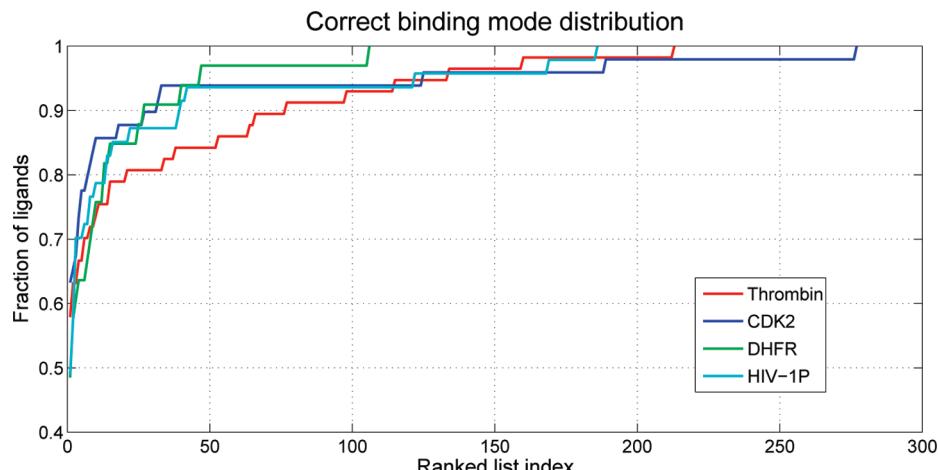


Figure 8. A cumulative plot of the indices at which ligands have their first predicted binding mode with RMSD < 2.5 Å. Ligands not having a correct binding mode within the top ranked 300 poses are excluded. If a ranked list is going to contain a correct binding mode within the top 300 poses it is likely to rank within the top 50.

Prediction with Decoys. We next considered the scenario where some ligands did not have any near-native binding modes in their ranked lists. We conducted a set of experiments similar to those of the previous section but allowed the input sets to include decoy ligands. We define a decoy ligand to be one whose list of docking algorithm generated binding modes does not contain a pose within 2.5 Å RMSD of the native binding mode. For all target proteins we use sets of 10 or 15 ligands having 1, 2, or 3 decoys and 50 binding modes in their ranked lists. Figure 7 illustrates the prediction results for these sets over 50 random experiments for each parameter setting (1200 experiments). Our method demonstrates average RMSD improvements over the naive approach of 1.05 Å, 0.79 Å, 0.92 Å, and 0.91 Å for the thrombin, CDK2, DHFR, and HIV-1P, respectively. These results illustrate the robustness of the method to the introduction of decoy ligands. In the tested cases, the decoys only moderately effect the accuracy of the algorithm.

Robustness of the Algorithm. The number of target systems for which we can evaluate algorithm performance is limited. The limitation is caused by the lack of experimentally determined native binding modes (i.e., the ground truth) for multiple ligands interacting with the same test protein. Therefore, we establish the performance of the algorithm by examining its performance over a range of modified input data sets. We evaluated the robustness of the algorithm with respect to the length of the ranked lists, ligand similarity, and sparsity of the search (i.e., the number of correct binding modes in the ranked lists).

Length of the Ranked Lists. We analyzed the virtual docking results and established the cutoff length for the ranked lists. Figure 8 shows the distribution of ligands with respect to the rank of their first correct binding mode (<2.5 Å). The docking results support the hypothesis that if a correct binding mode is generated by the docking algorithm that the correct pose is likely to appear within the first 50 ranked binding modes. We observed that ~95% of the time, if a ligand does not have a correct binding mode within its first 50, it will not have one ranked lower. This result is consistent with contemporary docking studies across multiple docking algorithms.^{5,9,11,31,32}

When the docking algorithm's top ranked pose is within 2.5 Å of native, it may still be possible to improve the quality

Table 2. Fraction of Ligands with a Correct Solution Ranked First (<2.5 Å) Which Also Have a Solution τ Percent Better within the Top 50 Ranked Candidate Binding Modes^a

protein	improvement threshold (τ)			
	>0%	10%	20%	30%
thrombin	0.88	0.76	0.70	0.55
CDK2	0.90	0.61	0.48	0.35
DHFR	0.94	0.69	0.56	0.31
HIV-1P	0.74	0.35	0.09	0.04
mean	0.87	0.60	0.45	0.31

^a See text. For 87% of the tested cases, when a correct binding mode appeared at the top of the ranked list, a pose more similar to the native binding mode appeared elsewhere in the top 50 poses.

of the predicted binding mode by searching lower in the ranked list of poses. We computed the frequency with which a pose could be identified with an RMSD τ percent better than the docking algorithm's top ranked pose. We define this cutoff percentage as the *improvement threshold* (Table 2). For the four systems tested, ligands having a correct binding mode ranked first are extremely likely (87%) to have a better binding mode further down the docking algorithm's ranked list. This result supports the overall motivation of our work. It suggests that in the majority of cases, the configuration of a set of ligands can be improved by selecting poses beyond the binding modes ranked highest by the underlying docking algorithm.

Ligand Similarity. We next investigated the impact of varying ligand similarity on the quality of the predictions. Similar ligands are likely to bind a target protein in similar fashion, thus they are likely to exploit similar interaction patterns¹⁹ and are conducive to the meta-analysis of our method. Renner et al. derived a similar conclusion in their binding mode prediction work, suggesting that the more similar the ligands, the better the binding mode prediction accuracy will be.²⁰ Therefore, there is a trade-off. Sets of extremely similar ligands provide little informational advantage over docking a single ligand. As the ligand set becomes more diverse, the independent docking runs of each ligand contribute different binding patterns to the meta-analysis and provide an informational advantage. However, in the limit, extremely diverse ligands may bind the receptor in completely different binding modes and may exhibit little self-

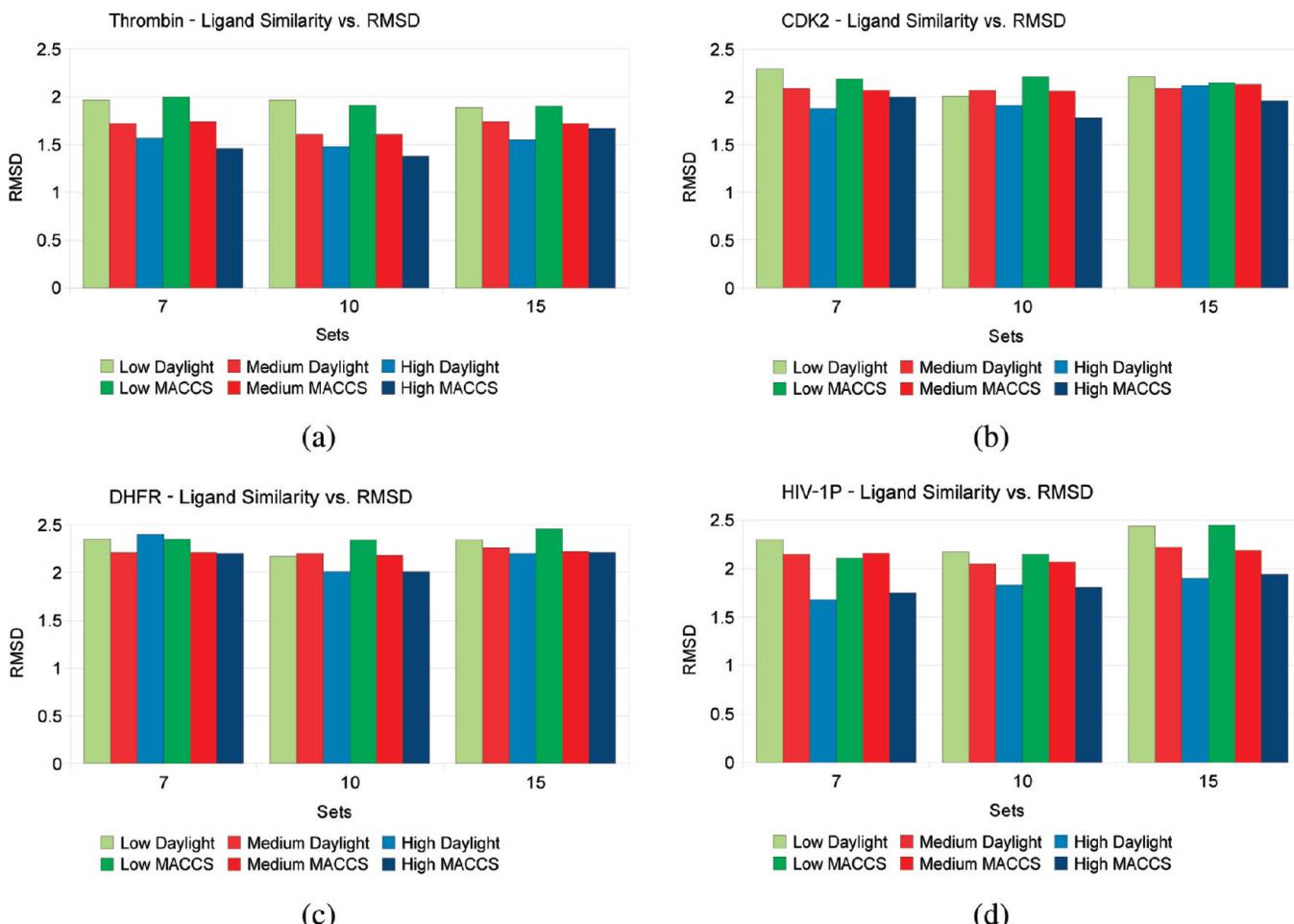


Figure 9. The changes in RMSD with respect to three bins of ligand similarity ranges over 1200 random experiments using ranked lists of length 30 and two different similarity fingerprints (Daylight and MACCS-like). There is a slight trend where ligand sets with higher similarity tend to generate configurations with lower RMSD.

consistency with respect to an autogenerated pharmacophoric map. Therefore, we studied the impact of ligand similarity on binding mode prediction.

We evaluated the ligand similarity and RMSD for 1200 randomly selected ligand sets. This evaluation accounted for the ability of our pharmacophore-based scoring approach to identify near-native configurations under varied ligand set similarities. We define the similarity of a set of ligands as the average pairwise similarity of ligands in that set. Pairwise ligand similarity is measured using the Tanimoto coefficient³³ of two types of molecular fingerprints implemented in the OpenBabel chemical toolbox.³⁴ The first is the Daylight fingerprints²⁶ which consists of linear segments of up to 7 atoms. The second is a set of SMARTS patterns defining chemical groups, similar to the MACCS fingerprints,³⁵ that represent the classification of organic compounds from the viewpoint of an organic chemist.

We measured the change in RMSD for ligands with varying degrees of similarity. The result of each randomly generated experiment was assigned to one of three bins reflecting *Low*, *Medium*, or *High* degrees of similarity. First, we calculated the mean similarity, μ , and standard deviation, σ , of all experiments with the same target protein, ligand set size, and fingerprint type. Then, the RMSD of the predicted configuration, with a ligand similarity value s was assigned to one of the three bins. Assignment was to the *Low* bin if $s < \mu - 1/\sigma$, *medium* bin if $\mu - 1/\sigma \leq s \leq \mu + 1/\sigma$, and *high* bin if $s > \mu + 1/\sigma$. Figure 9 summarizes the change in RMSD with respect to the ligand similarity bins. The results demonstrate a slight improvement in binding mode predictions for ligand sets with increased similarity. Nevertheless, the method is fairly robust to the changes in similarity showing an average RMSD improvement of 0.3 Å and 0.34 Å between the *Low* and the *High* ranges, for the Daylight and MACCS fingerprints, respectively. Interestingly, not only was the average improvement using the MACCS fingerprints higher than the Daylight ones but also the average Tanimoto coefficient values of the MACCS fingerprints were consistently higher than the Daylight ones (Table 3). This observation might be attributed to the way the Daylight and MACCS fingerprints account for ligand similarity. MACCS fingerprints account for the presence or

Table 3. Average Ligand Similarity of 1200 Randomly Generated Prediction Experiments^a

protein	average Tanimoto coefficients (std. dev.)	
	Daylight	MACCS-like
thrombin	0.38 (0.18)	0.47 (0.17)
CDK2	0.26 (0.13)	0.44 (0.14)
DHFR	0.43 (0.19)	0.62 (0.19)
HIV-1P	0.32 (0.23)	0.38 (0.21)

^a The table lists the average Tanimoto coefficient and standard deviation (std. dev.) of the ligand sets using Daylight and MACCS-like fingerprints.

+ $1/\sigma$, and high bin if $s > \mu + 1/\sigma$. Figure 9 summarizes the change in RMSD with respect to the ligand similarity bins. The results demonstrate a slight improvement in binding mode predictions for ligand sets with increased similarity. Nevertheless, the method is fairly robust to the changes in similarity showing an average RMSD improvement of 0.3 Å and 0.34 Å between the *Low* and the *High* ranges, for the Daylight and MACCS fingerprints, respectively. Interestingly, not only was the average improvement using the MACCS fingerprints higher than the Daylight ones but also the average Tanimoto coefficient values of the MACCS fingerprints were consistently higher than the Daylight ones (Table 3). This observation might be attributed to the way the Daylight and MACCS fingerprints account for ligand similarity. MACCS fingerprints account for the presence or

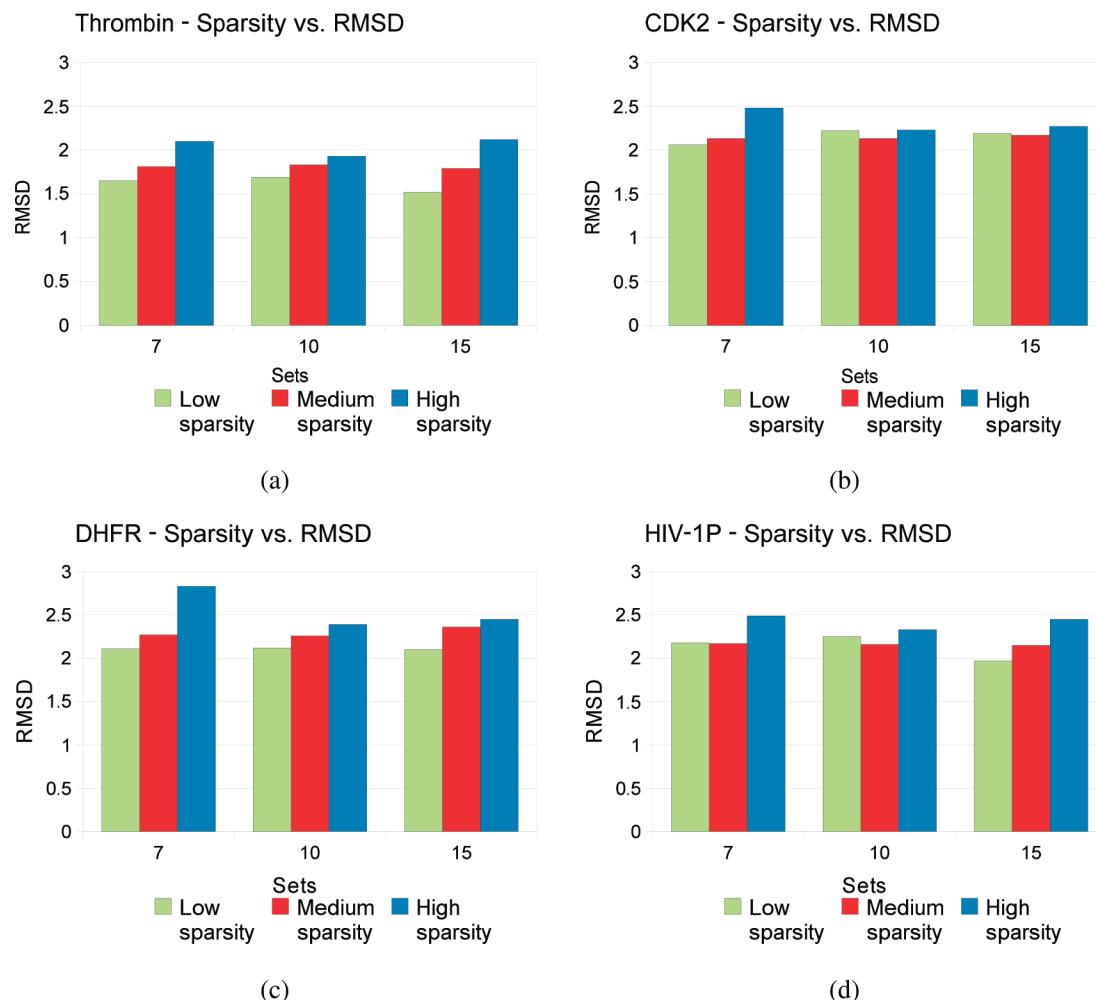


Figure 10. The changes in RMSD with respect to sparsity over 1200 random experiments using ranked lists of length 50. As expected, predicted configurations from pose sets with higher sparsity (i.e., fewer near-native poses) are less accurate than those with lower sparsity.

absence of similar chemical groups, whereas the Daylight fingerprints are more sensitive to structure. That is, structurally different ligands sharing similar sets of chemical groups will be similar under the MACCS fingerprints but less similar under the Daylight fingerprints. We suggest that the MACCS fingerprints better account for the degree of the expected coverage of pharmacophoric points by a ligand set. Diverse sets under the MACCS measurement imply diversity in the chemical groups, and diversity in the chemical groups implies more pharmacophoric points representing binding interactions. Thus, using the MACCS fingerprints, a similar set of ligands implies a better coverage of pharmacophoric points such that they become noticeable by the clustering phase.

Sparsity of the Search Space. The combinatorial optimization problem of identifying a good configuration becomes more difficult when the number of correctly identified binding modes per ligand is small. For example, consider a set of 7 ligands with 100 candidate binding modes where each ligand has only one correct binding mode (i.e., 99 poses have an RMSD > 2.5 Å). Such a set has only a single configuration (out of 10^{14} possible ones) for which all poses have RMSD < 2.5 Å.

We define the *sparsity* of a docking algorithm's ranked list as the ratio between the number of poses in the ranked list that are above a defined RMSD threshold (2.5 Å, in this study) and the total number of poses in that list. Additionally,

we define the sparsity of a set of ligands and docked poses to be the average sparsity over the ligands in that set. We measured the change in RMSD with respect to three bins of sparsity ranges (low, medium, high), in a method analogous to the ligand similarity analysis described before. Figure 10 summarizes the results of 1200 random experiments. When the sparsity of the ligand set increases (i.e., there are less correct poses in the ranked lists), the RMSD of the predicted configuration increases. The Low sparsity sets demonstrate an average improvement (over the High sparsity sets) of 21%, 7%, and 19% for sets of size 7, 10, and 15, respectively. This trend is maintained when evaluating the results for each protein system independently (data not shown). This phenomena may be explained by a trade-off between the coverage of pharmacophoric points and the size of the search space. A larger ligand set improves pharmacophoric coverage but increases the size of the search space. These results suggest that for the four systems tested, sets of 10 ligands demonstrate an optimal trade-off and thus provide the lowest sensitivity to the sparsity of the search.

CONCLUSION AND FUTURE WORK

We presented a novel algorithm to improve the binding mode predictions of a docking algorithm for the case of multiple ligands known to bind a common target. Our method

searches for a configuration (or set of binding modes, one per ligand) that is maximally self-consistent with its underlying pharmacophoric map. The consistency is scored using an objective function which considers the alignment of similar chemical groups. We showed that our objective function correlates well with the RMSD to the experimentally determined binding modes. Therefore, the configuration which optimizes our objective function should represent a set of binding poses maximally consistent with the true binding mode. The combinatorial optimization problem is addressed with a Gibbs sampling based technique which is initialized by Boltzmann weighted pose sampling. We performed two classes of experiments. In the first class, we considered the optimistic situation where each ligand was guaranteed to have at least one correct binding mode among the docking algorithm's list of poses. In the second class, we introduced decoys into the system. We then demonstrated the expected performance of the algorithm over different sets of input ligands varied by the size, similarity, and sparsity of their ranked lists. We demonstrated that across a range of conditions and four different protein systems, the configurations generated by our approach are 0.5–1.0 Å more accurate than the naive poses ranked highest by the underlying docking algorithm. An important advantage of our approach is its generality. While our method was demonstrated using the FlexX docking algorithm, our approach can utilize the ranked list of binding modes generated by any docking method. The results presented in this paper simply depend on the presence of a ranked list of binding modes, the majority of which should contain at least one near-native pose. We and others have demonstrated^{5,9,31} that most state-of-the-art docking algorithms satisfy this requirement. We emphasize that rather than aiming to replace any traditional docking algorithm, the goal of our work is to devise a meta-search technique capable of augmenting a docking algorithm. This generalization allows our method to be more easily incorporated into a range of drug discovery research pipelines.

Our method can be compared to meta-analysis methods, where information from multiple experiments is combined to improve performance. In the context of our work, the pharmacophoric map generated by a set of binding ligands allows the identification of correct binding modes, even when they are unfavorably ranked. The pharmacophore-based model discards poses which are favorably ranked by the docking algorithm (i.e., they have a high docking score) but which are not consistent with a common binding model. The approach is most successful when several similar yet distinct ligands bind a target protein using a common interaction profile. Because this phenomenon is consistent with previous docking and structural studies, the method may be generally applicable.

We present a step toward improving protein–ligand binding mode prediction for a set of ligands known to interact with a common protein. There is thus an important distinction between this work and traditional virtual screening algorithms. Whereas traditional approaches attempt to identify binding ligands from a large database of available compounds, our approach aims to more accurately predict the binding mode for a set of ligands which are already known to bind the target protein. We propose that our algorithm will be of use to the pharmaceutical industry in the structural analysis of molecular hits resulting from High Throughput

Screening (HTS) and Structure Activity Relationship (SAR) experiments.

There are a number of promising next steps for this work. We are currently considering augmenting the Gibbs sampling method with a more deterministic branch-and-bound search. While our current pharmacophoric model accounts only for the alignment of functional groups, a more elaborate model may integrate additional meta-data such as binding affinity. Furthermore, it may be possible to incorporate the negative binding results of an HTS experiment to eliminate pharmacophoric hypotheses and their corresponding configurations. We anticipate making a distribution of our method available in the near future.

ACKNOWLEDGMENT

We thank Mr. Abraham Heifets, Mr. Navdeep Jaitly, Mr. Satyam Merja, Ms. Maria Safi, and all members of the Lilien lab for helpful discussions and comments on drafts. This work is supported by a Bill and Melinda Gates Foundation (Grand Challenges Explorations) grant to R.H.L.

Note Added after ASAP Publication. This paper was published ASAP on August 27, 2009 with an error in the caption for Figure 4. The corrected version was published ASAP on September 1, 2009.

REFERENCES AND NOTES

- (1) Langer, T.; Hoffmann, R. D. *Pharmacophores and Pharmacophore Searches*, 1st ed.; Wiley-VCH: Palo Alto, 2006.
- (2) Patani, G. A.; LaVoie, E. J. Bioisosterism: A Rational Approach in Drug Design. *Chem. Rev.* **1996**, *96*, 3147–3176.
- (3) Congreve, M.; Murray, C. W.; Blundell, T. L. Keynote review: Structural biology and drug discovery. *Drug Discovery Today* **2005**, *10*, 895–907.
- (4) Bajorath, J. Integration of virtual and high-throughput screening. *Nat. Rev. Drug Discov.* **2002**, *1*, 882–894.
- (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) Sousa, S. F. F.; Fernandes, P. A. A.; Ramos, M. J. a. J. Protein-ligand docking: Current status and future challenges. *Proteins* **2006**, *65*, 15–26.
- (7) Kitchen, D. B.; Decornez, H.; Furr, J. R.; Bajorath, J. Docking and scoring in virtual screening for drug discovery: methods and applications. *Nat. Rev. Drug Discovery* **2004**, *3*, 935–949.
- (8) Chen, H.; Lyne, P. D.; Giordanetto, F.; Lovell, T.; Li, J. On Evaluating Molecular-Docking Methods for Pose Prediction and Enrichment Factors. *J. Chem. Inf. Model.* **2006**, *46*, 401–415.
- (9) Warren, G.; Andrews, C.; Capelli, A.-M.; Clarke, B.; LaLonde, J.; Lambert, M.; Lindvall, M.; Nevins, N.; Semus, S.; Senger, S.; Tedesco, G.; Wall, I.; Wolven, J.; Peishoff, C.; Head, M. A Critical Assessment of Docking Programs and Scoring Functions. *J. Med. Chem.* **2006**, *49*, 5912–5931.
- (10) Kontoyianni, M.; McClellan, L.; Sokol, G. Evaluation of Docking Performance: Comparative Data on Docking Algorithms. *J. Med. Chem.* **2004**, *47*, 558–565.
- (11) Cross, J. B.; Thompson, D. C.; Rai, B. K.; Baber, J. C.; Fan, K. Y.; Hu, Y.; Humbert, C. Comparison of Several Molecular Docking Programs: Pose Prediction and Virtual Screening Accuracy. *J. Chem. Inf. Model.* **2009**, *49*, 1455–1474, doi: 10.1021/ci900056c.
- (12) Barnum, D.; Greene, J.; Smellie, A.; Sprague, P. Identification of Common Functional Configurations Among Molecules. *J. Chem. Inf. Model.* **1996**, *36*, 563–571.
- (13) Richmond, N. J.; Abrams, C. A.; Wolohan, P. R. N.; Abrahamian, E.; Willett, P.; Clark, R. D. GALAHAD: 1. Pharmacophore identification by hypermolecular alignment of ligands in 3D. *J. Comput.-Aided Mol. Des.* **2006**, *V20*, 567–587.
- (14) Inbar, Y.; Schneidman-Duhovny, D.; Dror, O.; Nussinov, R.; Wolfson, H. J. Deterministic Pharmacophore Detection Via Multiple Flexible Alignment of Drug-Like Molecules. *Res. Comput. Mol. Biol.* **2007**, 412–429.
- (15) Joseph-McCarthy, D.; Thomas, B. E.; Belmarsh, M.; Moustakas, D.; Alvarez, J. C. Pharmacophore-based molecular docking to account for ligand flexibility. *Proteins* **2003**, *51*, 172–188.

- (16) Hindle, S. A.; Rarey, M.; Buning, C.; Lengauer, T. Flexible docking under pharmacophore type constraints. *J. Comput.-Aided Mol. Des.* **2002**, *16*, 129–149.
- (17) McGregor, M. A Pharmacophore Map of Small Molecule Protein Kinase Inhibitors. *J. Chem. Inf. Model.* **2007**, *47*, 2374–2382.
- (18) Chema, D.; Eren, D.; Yayon, A.; Goldblum, A.; Zaliani, A. Identifying the binding mode of a molecular scaffold. *J. Comput.-Aided Mol. Des.* **2004**, *18*, 23–40.
- (19) Bostrom, J.; Hogner, A.; Schmitt, S. Do Structurally Similar Ligands Bind in a Similar Fashion. *J. Med. Chem.* **2006**, *49*, 6716–6725.
- (20) Renner, S.; Derksen, S.; Radestock, S.; Morchen, F. Maximum Common Binding Modes (MCBM): Consensus Docking Scoring Using Multiple Ligand Information and Interaction Fingerprints. *J. Chem. Inf. Model.* **2008**, *48*, 319–332, doi: 10.1021/ci7003626.
- (21) Deng, Z.; Chuquui, 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.
- (22) Kahraman, A.; Morris, R. J.; Laskowski, R. A.; Thornton, J. M. Shape Variation in Protein Binding Pockets and their Ligands. *J. Mol. Biol.* **2007**, *368*, 283–301.
- (23) Wang, R.; Wang, S. How Does Consensus Scoring Work for Virtual Library Screening? An Idealized Computer Experiment. *J. Chem. Inf. Model.* **2001**, *41*, 1422–1426.
- (24) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E. The Protein Data Bank. *Nucleic Acids Res.* **2000**, *28*, 235–242.
- (25) Chen, X.; Rusinko, A.; Tropsha, A.; Young, S. S. Automated Pharmacophore Identification for Large Chemical Data Sets. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 887–896.
- (26) James, A. C.; Weininger, D.; Delany, J. *Daylight Theory Manual-Daylight 4.71*; 2000.
- (27) Schmitt, S.; Kuhn, D.; Klebe, G. A New Method to Detect Related Function Among Proteins Independent of Sequence and Fold Homology. *J. Mol. Biol.* **2002**, *323*, 387–406.
- (28) Frey, B. J.; Dueck, D. Clustering by passing Messages Between Data Points. *Science* **2007**, *315*, 972–976.
- (29) Bishop, C. M. *Pattern Recognition and Machine Learning (Information Science and Statistics)*; Springer: 2006.
- (30) Rarey, M.; Kramer, B.; Lengauer, T.; Klebe, G. A Fast Flexible Docking Method using an Incremental Construction Algorithm. *J. Mol. Biol.* **1996**, *261*, 470–489.
- (31) Onodera, K.; Satou, K.; Hirota, H. Evaluations of Molecular Docking Programs for Virtual Screening. *J. Chem. Inf. Model.* **2007**, *47*, 1609–1618.
- (32) Feher, M.; Williams, C. I. Effect of Input Differences on the Results of Docking Calculations. *J. Chem. Inf. Model.* 2009.
- (33) Tanimoto, T. T. *An Elementary Mathematical Theory of Classification and Prediction*; IBM Internal Report; New York, 1958.
- (34) Guha, R.; Howard, M.; Hutchison, G.; Murray-Rust, P.; Rzepa, H.; Steinbeck, C.; Wegner, J.; Willighagen, E. The Blue Obelisk—Interoperability in Chemical Informatics. *J. Chem. Inf. Model.* **2006**, *46*, 991–998.
- (35) Durant, J. L.; Leland, B. A.; Henry, D. R.; Nourse, J. G. Reoptimization of MDL Keys for Use in Drug Discovery. *J. Chem. Inf. Model.* **2002**, *42*, 1273–1280.

CI900199E