1	Incorporating structural similarity into a scoring function to enhance the
2	prediction of binding affinities
3	Beihong Ji, Xibing He, Yuzhao Zhang, Jingchen Zhai, Viet Hoang Man, Shuhan Liu, Junmei
4	Wang*
5 6	Department of Pharmaceutical Sciences and Computational Chemical Genomics Screening Center, School of Pharmacy, University of Pittsburgh, Pittsburgh, PA 15261, USA,
7	
8	* Corresponding author:
9	Junmei Wang, E-mail: junmei.wang@pitt.edu
10	
11	
12	
13	
14	
15	
16	
17	
18	
19	
20	
21 22	
14 15 16 17 18 19 20 21 22	

23 Abstract

24 In this study, we developed a novel algorithm to improve the screening performance of an arbitrary 25 docking scoring function by recalibrating the docking score of a query compound based on its 26 structure similarity with a set of training compounds, while the extra computational cost is 27 neglectable. Two popular docking methods, Glide and AutoDock Vina were adopted as the 28 original scoring functions to be processed with our new algorithm and similar improvement 29 performance was achieved. Predicted binding affinities were compared against experimental data 30 from ChEMBL and DUD-E databases. 11 representative drug receptors from diverse drug target 31 categories were applied to evaluate the hybrid scoring function. The effects of four different 32 fingerprints (FP2, FP3, FP4, and MACCS) and the four different compound similarity effect (CSE) functions were explored. Encouragingly, the screening performance was significantly improved 33 for all 11 drug targets especially when $CSE=S^4$ (S is the Tanimoto structural similarity) and FP2 34 35 fingerprint were applied. The average predictive index (PI) values increased from 0.34 to 0.66 and 36 0.39 to 0.71 for the Glide and AutoDock vina scoring functions, respectively. To evaluate the 37 performance of the calibration algorithm in drug lead identification, we also imposed an upper 38 limit on the structural similarity to mimic the real scenario of screening diverse libraries for which 39 query ligands are general-purpose screening compounds and they are not necessarily structurally 40 similar to reference ligands. Encouragingly, we found our hybrid scoring function still 41 outperformed the original docking scoring function. The hybrid scoring function was further 42 evaluated using external datasets for two systems and we found the PI values increased from 0.24 43 to 0.46 and 0.14 to 0.42 for A2AR and CFX systems, respectively. In a conclusion, our calibration 44 algorithm can significantly improve the virtual screening performance in both drug lead 45 optimization and identification phases with neglectable computational cost.

Keywords: Scoring Function; Virtual Screening; Docking; Compound Similarity; Fingerprint;
 Protein-Ligand binding; Computer-Aided Drug Design

48

49 Introduction

50 In order to save time and cost in drug discovery projects, various in silico approaches have 51 been developed and applied to reduce the number of compounds which are to be experimentally 52 synthesized and tested. Among these computer-aided drug design/discovery (CADD) methods, 53 virtual screenings of large chemical databases for potential bioactive molecules are usually 54 conducted at the very beginning stage to rapidly narrow down the candidates from millions of 55 compounds to manageable numbers (thousands or hundreds). Depending on whether the 3D 56 structural information of the target receptor is available and utilized, virtual screening (VS) 57 methods can be broadly classified into structure-based (SBVS) and ligand-based (LBVS) [1]. 58 Docking & scoring is a typical SBVS method, which predicts whether a ligand can favorably 59 interact with a receptor (protein or nucleic acid) at its binding site, and if yes, the binding mode 60 and binding affinity measured by the docking scoring function are determined. [1, 2]. Similarity 61 search is a typical LBVS method, which predicts activity of query compounds depending on their 62 similarities/dissimilarities to known reference ligands by utilizing numerical similarity descriptors 63 (fingerprints) [3]. Both docking and similarity methods have been successfully carried out 64 independently or hierarchically to screen out confidently inactive compounds for specific receptors 65 of interest. Compared to docking, similarity search is even faster, and therefore it is suitable to 66 filter super large database before docking takes place. However, both similarity search and docking 67 suffer from poor accuracy to rank potentially active ligands and prioritize top candidates to be 68 suggested for experiments. Similarity search is based on the Similarity Property Principle (SPP) 69 [4], i.e., structurally similar molecules are likely to possess similar biological properties and

70 activities. However, this hypothesis is not always true. Sometimes small chemical differences may 71 arise highly different activity ('activity cliffs') [5]. Accuracy of docking methods is limited due to 72 lack of modelling structural flexibility of target receptors, effects of solvation and entropy changes, 73 etc. These limitations of docking & scoring methods may be overcome by more accurate 74 methodologies, such as end-point methods (MM-PBSA, MM-GBSA, LIE, etc.) [6, 7], or rigorous 75 alchemical free energy methods (FEP, TI, etc.) [8, 9], with the price of much higher computational 76 cost and much longer time. But we cannot help wondering: is there a way to improve the accuracy 77 of docking & scoring methods without much more extra computational cost? With the advances 78 of high-throughput screening technique, more and more compounds were measured against a drug 79 target. ChEMBL [10] is a curated database which collects binding affinities of bioactive molecules 80 for a drug target. How can we utilize the information on the known structures and activities to 81 improve screening performance? Secondly, can we combine docking & scoring methods with the 82 extremely fast methods of similarity calculations to improve the accuracy of binding affinity 83 estimation? If so, how can we incorporate the two types of scores into one hybrid scoring function? 84 In this work, we attempted to develop a novel algorithm to make a good use of those valuable 85 information on known bioactive compounds.

Docking programs utilize scoring functions to estimate the binding affinities. Currently, scoring functions can be generally classified into following four categories based on how proteinligand energy is predicted: (1) force-field-based, (2) empirical, (3) knowledge-based, (4) descriptor-based [11]. Although the underlying mechanisms of four categories are different, all of those developed scoring functions are trying to pursue promising results in the protein-ligand binding prediction. Considering the progress and achievements that have already been made by the existing scoring functions during the past few decades, we believe it will be more valuable to

93 make some improvements based on current developed scoring functions. Therefore, in this work, 94 instead of developing a completely novel scoring function for the binding affinity prediction, we 95 introduced a more practical and universal approach that can improve the scoring power for an 96 arbitrary docking scoring function. Our new scoring algorithm is suitable to the following scenario: 97 the binding affinities against a specific receptor for a certain number of compounds (hereafter 98 called reference compounds or reference ligands) have been experimentally measured, but many 99 more compounds (hereafter called query compounds) need to be estimated in silico. Such a 100 scenario is frequently encountered in real drug design projects. Two components were 101 incorporated to compose our new algorithm: (1) the structural difference between the query 102 compounds and reference ligands; (2) the deviation of the docking scores of those reference ligands 103 and their experimental binding affinities. Both the above components can serve as weights to 104 calibrate the original docking scores of ligands.

105 Materials and Methods

106 Work outline. The outline of our work is shown in the flow chart of Fig. 1. To evaluate the 107 feasibility and performance of our algorithm, we collected the X-ray crystal structures of 11 108 receptors of various categories from the Protein Data Bank [12] (https://www.rcsb.org) and their 109 available ligand data from the ChEMBL database [13, 14] (https://www.ebi.ac.uk/chembl/). 110 Collected compounds for each receptor were divided into reference set (reference ligands) and 111 validation set (query compounds). See Table 1 and Supplementary Information. Compounds were 112 docked to their corresponding targets with a state-of-art docking program Glide, which was 113 selected as the basic scoring function for further improvement [15, 16]. The Tanimoto Coefficient 114 Tc [3] was calculated by utilizing the Open Babel program version 2.3.1 (http://openbabel.org) [17]. Four popular 2D fingerprints (FP2, FP3, FP4, MACCS) [17] available in the Open Babel 115

116 package were adopted and compared in this study. The proposed hybrid scoring function was 117 applied to calibrate the Glide docking scores. The scoring and ranking power of the hybrid scoring 118 function as well as the original Glide docking scoring function was measured by root-mean-squareerror (RMSE), mean-absolute-error (MAE), Pearson's correlation coefficient (R²), and the 119 120 predictive index (PI) between the docking scores and the experimental binding affinities [18-20]. 121 In addition, we evaluated the screening power of the hybrid scoring function by using enrichment 122 factor (EF) and the area under the curve (AUC) of receiver operating characteristic (ROC) curve 123 [21]. We also explored the choice of fingerprint type and CSE function form to optimize the 124 screening performance. More details about the preparation of receptors and ligand datasets, 125 docking software and procedures, and methods of evaluations are described in the following 126 sessions.

127 **Preparation of receptor datasets.** To test the developed algorithm in this study, 11 receptors 128 were selected according to the experimental records in ChEMBL database. The receptors can be 129 divided into three classes: (1) 6 top receptors from diverse categories; (2) 4 top receptors from 130 GPCR family; (3) one RNA receptor. Coagulation factor X (CFX), dopamine D2 receptor (D2R), 131 μ opioid receptor (MOR), Extracellular signal-regulated kinase 2 (ERK2), vascular endothelial 132 growth factor receptor 2 (VEGFR2) and estrogen receptor (ER) are in class (1). Among them, CFX 133 is a member of Protease. D2R and MOR are from family A GPCR. EKR2 and VEGFR2 are 134 members of the kinase family, while ER is the nuclear receptor. In class (2), there are serotonin 135 2A receptor (5HT2AR), adenosine A2A receptor (A2AR), cannabinoid receptor 1 (CB1) and 136 muscarinic acetylcholine M1 receptor (M1R). Finally, in class (3), the ribosomal RNA (rRNA) A-137 site was selected as the receptor. All the receptor categories were selected according to the rank of 138 member amounts in their category families, while the receptors themselves were selected based on

the compound number recorded in binding assays. The above information was collected from
ChEMBL database and is shown in Table 1. Then the X-ray crystal structures of the above 11
target receptors were retrieved from Protein Data Bank (detail information was shown in Table
S1). The sources of all targets are Homo sapiens.

143 **Preparation of ligand datasets.** To better compare the binding energy of each ligand that 144 binds to the same receptor, we collected 3D structures (SDF format) of ligands with the inhibition 145 constant (K_i) values recorded from binding assays in the ChEMBL database. For the ribosomal 146 RNA receptor, because of the limited number of K_i activities, compounds with dissociation 147 constant (K_d) values were collected. It is of note that the activities of those collected compounds 148 for each receptor were measured using the same methods. To balance the distribution of K_i 149 activities for each receptor, compounds were hierarchically classified into 4 levels according to 150 their K_i values: K_i < 10 nM, 10 nM \leq K_i < 1 μ M, 1 μ M \leq K_i < 100 μ M and K_i \geq 100 μ M. In 151 each level, 300 or less than 300 compounds (if the number of compounds in the level does not 152 reach 300) were randomly collected by utilizing numpy.random.choice in Python 3.7 program 153 [22]. For one selected compound with 2 or more K_i values that came from various assays, the 154 average K_i was used. To evaluate the screening power of our approach, we categorized the selected 155 compounds into the active and inactive sets by the cutoff of $K_i/K_d = 100$ nM. This value is lower 156 than normal threshold, 10 μ M, but it can better balance the numbers of compounds in the active 157 and inactive sets. The experimental K_i/K_d for each collected compound was then converted to the 158 experimental ligand-receptor binding energy (kcal/mol) by the Equation 1.

159
$$\Delta G_{binding} = -RT ln K_{i(d)} (1)$$

160 where $\Delta G_{binding}$ is the binding energy of the ligand, R is the gas constant with a value of 8.314 161 $J \cdot mol^{-1} \cdot K^{-1}$ and T is the room temperature under standard pressure with the value of 298.15 162 K.

163 To exclude compounds with very weak binding affinities and make our evaluation more 164 reliable, compounds for 11 receptors with experimental binding energy higher than -4 kcal/mol 165 were removed from the selected datasets. Then we randomly separated compounds into training 166 datasets and testing datasets by the proportion of 4:1. The number of compounds in training and 167 testing sets for each receptor was listed in Table 1. The structures of compounds were not only 168 stored in the mol2 format but also converted into different 2D fingerprints for the further 169 exploration of our algorithm by Open Babel program, which is an expert chemical toolbox for the 170 format interconversion of chemical data [23]. Four 2D fingerprints are available in OpenBabel 171 according to its documentation (http://openbabel.org/docs/dev/Features/Fingerprints.html): (1) 172 FP2, a path-based fingerprint stored in a 1024-bit vector; (2) FP3, s series of SMARTS queries 173 stored in 55 bits; (3) FP4, s series of SMARTS queries stored 307 bits; (4) MACCS, a series of 174 SMARTS patterns stored in 166 bits.

To critically evaluate the performance of our calibration algorithm, we performed extra validation test on hundreds of compounds with activities (K_i) of A2AR and CFX from an additional database, DUD-E database [24, 25]. After excluding the compounds with binding energy higher than -4 kcal/mol and those already included in the reference datasets, two external test datasets which have 1973 and 1599 unique compounds were compiled for the A2AR and CFX systems, respectively. All compounds from DUD-E were treated as query molecules and their docking scores were calibrated with the compounds in the corresponding ChEMBL dataset as references.

182 **Docking software and procedure.** We docked selected compounds (include training sets, test 183 sets and external test sets) to their corresponding receptors utilizing the Glide docking program 184 implemented in the Schrodinger software (Maestro 11.2). Before docking, the downloaded SDF 185 files of ligands were processed with the LigPrep module in Maestro. The downloaded PDB files 186 of receptors were processed with the module of *Protein Preparation Wizard* in Maestro: removing 187 co-crystallized solvent and ions, adding hydrogen atoms and missing site-chain atoms, energy 188 minimization on hydrogen atoms. Then we defined the binding site based on the geometric center 189 of the native bound ligand without taking constraint or rotatable group into consideration. The 190 flexible docking with post-docking minimization for 11 systems was conducted by the following 191 settings: van der Waal radius scaling factor was 0.80, the partial charge cutoff for ligands was 0.15, 192 the intramolecular hydrogen bond formation was rewarded, the number of poses per ligand to 193 include was 10. The top binding pose with the best docking energy score was retained and stored. 194 To investigate the impact of different docking scoring function on our calibration algorithm, 195 the performance of our hybrid scoring functions was also evaluated for the AutoDock Vina 196 docking scoring function [26]. Again, selected compounds were docked to their corresponding 197 receptors utilizing the AutoDock Vina docking program. The receptor preparation was performed 198 following the same protocol in Glide docking program. The binding site and space were defined 199 based on the geometric center and the size of the native bound ligand without taking constraint or 200 rotatable group into consideration. Considering the different docking mechanisms of two docking 201 programs, this time, compounds with experimental binding energy higher than -5 kcal/mol were 202 removed from the selected datasets to exclude compounds with weak binding affinities. Then 203 compounds were randomly separated into training datasets and testing datasets proportionally and 204 the docking score was calibrated using our proposed calibration approach. The Glide docking

scores and AutoDock Vina docking scores as well as the experimental binding affinities (converted
from K_i values) of compounds in the reference and validation sets were listed in Table S2.

Algorithm for docking score calibration. The new algorithm we proposed to calibrate the docking scores from a normal docking program is described as below:

209
$$DS_j = DS_j^0 \left[\frac{1}{\omega} \sum_{i \neq j}^n S_{ij}^p \frac{\Delta G_i}{DS_i} \right]$$
(2)

210
$$\omega = \sum_{i \neq j}^{n} S_{ij}^{p} \quad (3)$$

where DS_j^0 and DS_j are the docking score of the j^{th} query compound before and after the 211 calibration. S_{ij} is the structural similarity between the j^{th} query compound and the i^{th} reference 212 ligand. The exponent p is treated as an integer constant with its value varying from 1 to 4 in this 213 study, for the exploration of the developed formula. We referred S_{ij}^p as compound similarity effect 214 (CSE) function for convenience of discussion. n is the total number of reference ligands in the 215 reference dataset. DS_i is the docking score of the i^{th} reference ligand. ΔG_i is the experimental 216 binding energy (kcal/mol) of the *i*th compound in the reference dataset, which is converted from 217 218 the experimental K_i/K_d by the Equation (1).

In this study, the structural similarity S_{ij} between two compounds is represented by the Tanimoto Coefficient (Tc) calculated from their 2D fingerprints: [3]

221
$$T_c(X,Y) = \frac{z}{x+y-z}$$
 (4)

Where x and y are the number of bits in the fingerprints of compounds X and Y, z is the number of bits set shared by compounds X and Y. Tc has a range between 0 to 1 and a larger value means higher structural similarity between two compounds. The Tc calculation was carried out by utilizing Open Babel under a Python 3.7 environment. Given a simple example to demonstrate how the algorithm works, we assume there are only two reference compounds, i and j, whose docking scores are -8.0 and -10.0 kcal/mol and their experimental values are -7.0 and -8.0 kcal/mol, respectively. The docking score of the query compound is -9.0 and the similarity between the query compound and reference compound i and jare 0.9 and 0.5, respectively. Assuming p is 4, then after the calibration, the new docking score for the query compound becomes:

232
$$DS_{cali} = -9.0 \frac{1}{0.9^4 + 0.5^4} \left[0.9^4 \frac{-7.0}{-8.0} + 0.5^4 \frac{-8.0}{-10.0} \right] = -7.82$$

Apparently, reference compound *i* has more impact than *j* in the docking score calibration for this query compound. It is noted that our algorithm may not always improve the performance of docking score. There is a possibility that the docking score becomes worse after the calibration. However, we expect the similarity-based calibration can improve the binding affinity prediction in most scenarios with a certain size of the reference set.

238 **Performance evaluation.** To reduce the systematic error during the calculation, the random 239 separation of compounds into training and testing sets before the calibration was repeated for ten 240 times for each target. The mean value and 95% CI for all performance metrics were then calculated. 241 To evaluate the scoring and ranking performance of our algorithm, for each receptor, the docking 242 score of compounds in the test set was compared to their experimental energies individually utilizing four different measurements, RMSE, MAE, R² and PI. By comparing the mean calibrated 243 244 docking score with the original docking score, the scoring function was considered to be improved if RMSE and MAE reduced, while R² and PI increased. We also calculated the difference between 245 the calibrated docking score and the original one. The difference is respectively represented by 246 247 dRMSE, dMAE, dR^2 and dPI. For the evaluation of screening power, the area under the curve (AUC) of receiver operating characteristic (ROC) curve and enrichment factor (EF) at 10% 248

249 $(EF_{10\%})$ and 40% $(EF_{40\%})$ levels were adopted as the performance metrics. In comparison to the 250 simple docking scoring function, our new calibration algorithm on docking score was considered 251 to have better screening power if AUC and EF increased. We applied the same protocols to 252 evaluate the performance of the calibration algorithm on the datasets of A2AR and CFX targets 253 from the DUD-E database, except that the enrichment factors were calculated with different hit 254 rates. Considering the sample sizes of datasets for these two targets are relatively large in the DUE-255 E database, we utilized $EF_{1\%}$ and $EF_{10\%}$ to evaluate the screening power. The equations for the 256 calculation of metrics PI and EF were described in detail in Supplementary Information.

257 We defined two scenarios, "focus library" and "diverse library", which are respectively 258 appliable to drug lead optimization and lead identification in drug discovery, to evaluate the 259 algorithm by limiting the range of Tc. In other words, the training compounds which did not meet 260 the criterion of Tc range were excluded from calculations for the calibration. In the "focus library" 261 scenario, for each system, we set up a lower bound Tc value. Below this threshold, Tc will be too 262 low to improve the performance of the scoring function. In the "diverse library" scenario, besides 263 the lower bound Tc value, we also set up an upper bound for Tc. We randomly collected 10,000 264 screening compounds from ZINC database [27] (https://zinc.docking.org/) and calculated the 265 structural similarity between testing compounds and screening compounds individually. The Tc 266 value at which more than a half of screening compounds are lower than is selected as the upper 267 bound Tc value. It is noted that this is a very stringent method to determine the Tc upper bound 268 value.

269 **Results and Discussion**

The impact of fingerprint and CSE function. We first studied the calibration performance
using the Glide scoring function. For all 11 systems, the scoring power measured by RMSE and

MAE and ranking power measured by R^2 and PI of the original docking scoring function and the 272 hybrid scoring functions applying different fingerprints and CSE function are shown in Tables 273 274 S3-S4 and Fig. 2. According to Tables S3-S4, the developed algorithm can improve the accuracy 275 of original docking score for most of systems, no matter what fingerprint was used or what CSE 276 function was adopted. Specifically, when FP2 fingerprint was used for the similarity calculation 277 between compounds in our algorithm, the docking scores can improve, regardless of types of systems or CSE function. Similarly, when $CSE = S_{ij}^4$, the performance of the scoring function 278 279 enhanced for all 4 fingerprints. The comparison among the performance of the algorithms when 280 employing different fingerprints and CSE functions is clearly illustrated in Fig. 2. For most 281 systems, the enhanced effect of the algorithms with different calibrating functions can be ranked 282 from the largest to the smallest when the p value varied from 4 to 1. As to fingerprint type, Fig. 2 also demonstrated that FP2 stands out as it mostly has lower RMSE and MAE, higher R² and PI 283 284 than other fingerprints.

To test the generalizability of our calibration algorithm, we also studied the calibration performance using docking scores generated by another commonly-used docking program, AutoDock Vina. As shown in **Table S5** and **Fig. S1**, the same conclusion was reached for this docking program, i.e., the S⁴ outperforms other CES functions and FP2 outperforms other fingerprint types.

It is easy to understand why the performance of the algorithm became better when p value increased. As p value rises, the impact of the reference compounds that are structurally similar to the query compounds increases, meanwhile the impact of the reference compounds with low similarity reduces. As such, the weight of the similarity contribution will boost if the power of the similarity increases in the formula. It can be expected if the p value is higher than 4, the algorithm 295 might continue improving the scoring function even in a more positive way. However, to balance 296 the contributions from both the original docking scores and compound similarity effect, we let the 297 maximal p value stop at 4.

298 Another interesting factor that can affect the performance of the algorithm is the type of 299 fingerprints. The different underlying mechanisms of those fingerprints are likely to explain their 300 different effects. Unlike FP3, FP4 and MACCS that are substructure-based fingerprints based on 301 sets of SMARTS patterns, FP2 is a path-based fingerprint that indexes small molecules fragments 302 based on linear segments of up to 7 atoms, which might elucidate why FP2 performed better in our 303 algorithm [28]. As FP2 is more specific and can be used in any initial chemical searches, we 304 assume the similarity calculation based on FP2 is able to amplify the weights of those structurally 305 similar references to a greater extent and better offset the shortage of traditional scoring function. 306 For example, the traditional docking score is always averagely high even for those ligands with 307 relatively low binding affinity, leading to insufficient differentiation of docking results. On the 308 other hand, the different performances of those three substructure-based fingerprints (FP3, FP4 309 and MACCS), are likely caused by more complicated reasons. One probable reason is their 310 different number of descriptors. For example, utilizing FP3 improved the scoring function very 311 limitedly, which might be explained by its limited number of bits of 55 versus FP4 with 307 bits 312 and MACCS with 166 bits stored in Open Babel [29].

The calibration performance using the best hybrid function is similar for the two docking scoring functions. As shown in **Table 2**, for the tested 11 receptors with related thousands of ligands, on average, MAE decreased from 2.05 (Glide) to 1.34, 1.74 (AutoDock Vina) to 1.15; RMSE decreased from 2.53 (Glide) to 1.69, 2.11 (AutoDock Vina) to 1.47; while R² increased from 0.14 (Glide) to 0.44, and from 0.16 (AutoDock Vina) to 0.50; PI increased from 0.34 (Glide) to 0.66, and from 0.39 (AutoDock Vina) to 0.71. In other words, the improvement for average values of mean MAE, RMSE, R^2 and PI over 11 receptors are 0.71 kcal/mol, 0.84 kcal/mol, 0.30, and 0.32 respectively for the Glide docking scoring function, and the corresponding values are 0.59 kcal/mol, 0.64 kcal/mol, 0.34 and 0.32 for the AutoDock Vina scoring function. Interestingly, the boost performance on scoring power (RMSE and MAE) is better for Glide docking scoring function, while the ranking power (R^2 and PI) is better for AutoDock Vina.

324 On the other hand, the improvement of performance on screening power by using our 325 calibration algorithm further validated our approach. As shown in Table 3, EF_{10%} and EF_{40%} of 326 the docking results after the calibration are better than the results before the calibration for all the scenarios except for $EF_{10\%}$ of VEGFR2, for which the value before the calibration is slightly better 327 328 (2.18 vs 2.15). Considering the small sample size for rRNA receptor, it is reasonable to find 329 significant larger enrichment factor values in the calculated metrics. After excluding the rRNA 330 target, on average the mean $EF_{10\%}$ and $EF_{40\%}$ increased about 25% and 22% after the calibration 331 of docking scores. Similarly, mean AUC of the docking results after the calibration is significantly 332 improved for all the targets, as shown in Table 3 and illustrated in Fig. 4. Without considering the 333 performance on rRNA target, on average the AUC improved approximately 20% for the docking 334 results after the calibration. Of note that the measurement of screening power may be biased for 335 some drug target due to the imbalance between the number of actives and inactives, such as the 336 actives only account for about 12% of the total compounds in the test sets.

The impact from receptor categories on the calibration performance. As shown in Table and Figs. 2 and S1, the basic performance of docking score for all systems is different. Take the Glide docking scoring function as an example, for ERK2 drug target, the performance of original docking results is acceptable with a low mean MAE (1.20 kcal/mol) and RMSE (1.55 kcal/mol),

and a relatively high mean R^2 (0.35) and PI (0.74). On the other hand, for some receptors such as 341 342 MOR, the performance of the original docking results is not satisfying with a high mean MAE 343 (3.28 kcal/mol) and RMSE (3.94 kcal/mol), and a low mean R² (0.02) and PI (0.11). Therefore, in 344 order to further compare the improved effect of the algorithm between different systems by 345 excluding the impact from initial baselines of various systems, we quantitatively estimated the 346 difference between the calibrated docking score and original docking score using parameters dRMSE, dMAE, dR² and dPI which quantitively measure the difference of those measurements 347 348 before and after the calibration (Table S4 and Fig. 3). It is observed that the extent of the 349 improvement by using this algorithm varied from system to system, which indicates that although 350 the developed algorithm can be adaptable for various receptors, its enhanced effect can still differ 351 and depends on the receptor to some extent. The improvement on the scoring power and ranking 352 power is more prominently for those systems with poor docking performance, such as D2R, MOR, 353 5HT2AR, and CB1. After the calibration using the best hybrid scoring function, the PI increases 354 by 356%, 273%, 174% and 220% for the four systems correspondingly (Table 2). After the 355 calibration, not only the docking performance is enhanced, but also the standard deviations and 356 the CIs of the metrics measuring the docking performance are decreased among different drug 357 receptors.

Application of calibration in drug lead identification. In above, we discussed our hybrid scoring function can enhance screening performance for focused compound libraries in drug lead optimization. Next, we evaluated how well the best hybrid function (FP2 fingerprint and CSE=S⁴) performs for diverse compound libraries in drug lead identification using two studies. First, we created "diverse libraries" by imposing an upper limit of Tc as described in Methods section. When we calibrated the docking score of a query compound, we applied an upper limit of Tc value to exclude reference compounds which are structurally similar to the query compound to participate calibration. Encouragingly, even applying relatively low upper bound Tc values, our best hybrid scoring function can still enhance the docking performance for all receptors as shown in **Table 4**, even though the extend of the enhancement becomes much smaller as expected. The average values of mean MAE and RMSE decreased by 14.1%, 14.7%, respectively; and R² and PI increased by 26.7% and 17.1%, respectively, after the calibration. In a real situation, the enhancement may be significantly larger as demonstrated in the second study.

371 In the second study, we recalculated docking scores of the Glide scoring function for a set of 372 external test compound libraries collected by DUDE-E database for two drug targets, A2AR and 373 CFX. Unlike the first study, we did not impose an upper bound of Tc in selecting reference 374 compounds to mimic the real situation in virtual screening studies, however, for the test compound 375 libraries, we eliminated all the entries which were duplicated with reference compounds. The 376 performance of the best hybrid docking scoring function is summarized in Tables S6-S7 and 377 shown in Figs. S2-S3. The MAE and RMSE were respectively dropped from 1.44 to 1.05, and 378 1.78 to 1.37 kcal/mol for A2AR; and the two scoring power metrics were decreased from 2.71 to 1.66 and 3.25 to 2.13 kcal/mol for CFX. Similarly, the ranking power metrics R² and PI were also 379 significantly increased for both systems. For A2AR, R² changed from 0.05 to 0.19 and PI changed 380 381 from 0.24 to 0.46 (a 92% increase); and for CFX, R² changed from 0.02 to 0.16, and PI changed 382 from 0.14 to 0.42 (a 200% increase). As for the screening power, the EF_{1%} and EF_{10%} respectively 383 enhanced from 1.18 to 1.35 and from 1.12 to 1.32 for CFX, while these two metrics 384 correspondingly increased from 0.98 to 1.06 and 1.13 to 1.28 for A2AR. Last, the AUC values 385 were increased from 0.58 to 0.71 for CFX and from 0.61 to 0.71 for A2AR.

386 Taken together, in the scenario of drug lead identification, our calibration algorithm can still 387 significantly improve the docking performance measured by MAE and RMSE for scoring power, R² and PI for the ranking power and EF and AUC for screening power. On the other hand, we 388 389 pointed out that our method is based on docking results, hence the final performance on ranking 390 compounds after our calibration algorithm may not meet the high standards of correctly ranking 391 and prioritizing top compounds in the next stage of lead optimization, for which the more rigorous 392 but much more expensive methods, such as alchemical free energy calculation using free energy 393 perturbation [30] and thermodynamic integration [31, 32], are usually adopted.

394

395 Conclusion

396 In summary, we developed a novel algorithm for quickly improving the scoring power and 397 ranking power of a general scoring function used in a docking program by calibrating the docking 398 score according to the structural similarities between the query compound and a set of reference 399 compounds, whose experimental binding affinities have been measured. To evaluate the 400 performance of the algorithm, we collected 11 receptors from different categories and estimated 401 the enhanced effect of our algorithm on the Glide docking score by utilizing merit metrics of RMSE, MAE, R², PI, EF and AUC. Fingerprint type and structural similarity effect function, the 402 403 two factors which can significantly impact the performance of the calibration algorithm were 404 systematically explored. The results showed that our algorithm could enhance the performance of 405 the original docking scoring function in both the focused-library and diverse-library scenarios. We found that a combination of using FP2 fingerprint and a S⁴ CSE function can maximize the 406 407 calibration performance for most systems. For the scenario of using the hybrid scoring function in 408 drug lead optimization, the PI increased by 0.32 for both the Glide and AutoDock Vina scoring

409 functions; for the scenario of using the hybrid scoring function in drug lead identification, the PI 410 values of docking screenings using external test sets increased by 0.22 and 0.28 for A2AR and CFX systems, respectively. Thus, we successfully developed an algorithm which integrates 411 412 structure-based docking scores and ligand-based structural similarity scores into a hybrid scoring 413 function and make a good use of known experimental values. With more and more measured 414 binding affinity data collected by public databases like ChEMBL, our calibration algorithm could 415 have more and more broad applications in structure-based drug design. Afterall, the significantly 416 enhanced performance is achieved by a simple calibration algorithm whose computational cost is 417 neglectable.

418

419 Supplementary Information

420 Supplementary Information. Table S1 lists the name, entry code, resolution, released date and 421 deposition author for each receptor studied in this paper. Table S2 lists the Glide docking scores 422 (Table S2A) and AutoDock Vina docking scores (Table S2B) as well as the experimental binding 423 affinities of compounds in the reference and validation sets. Table S3 lists the RMSE, MAE, R² 424 and PI values before and after calibration of the Glide docking scores under the conditions of 425 different CSE function and fingerprint. Table S4 lists the difference of metrics for the 426 measurement of docking performance before and after the calibration, i.e., dRMSE, dMAE, dR² 427 and dPI for the Glide scoring function. Table S5 lists and Fig. S1 shows the RMSE, MAE, R² and 428 PI values before and after calibration of the AutoDock Vina docking scores under the conditions 429 of different CSE functions and fingerprints. Table S6 shows RMSE, MAE, R² and PI values before 430 and after calibration of Glide docking scores for compounds in the external test sets from DUD-E 431 database. Fig. S2 shows the comparison of RMSE, MAE, R2 and PI values before and after the 432 calibration of the Glide docking scores for A2AR and CFX external test sets using the best hybrid

- 433 scoring function (FP2 fingerprint with CSE=S4). Table S7 displays AUC, EF_{1%} and EF_{10%} values
- 434 before and after calibration of Glide docking scores for compounds in the external test sets from
- 435 DUD-E database. Fig. S3 shows ROC curves before and after calibration of the Glide docking
- 436 scores for A2AR and CFX external test sets.
- 437

438 Availability of data and materials

- 439 All data come from publication domain. The associated parameters of the hybrid scoring
- 440 functions for each drug target were presented in tables.

441 **Competing Interests**

442 The authors declare no competing financial interest.

443 Funding

- 444 This work was supported by the funds from National Science Foundation (1955260) and National
- 445 Institutes of Health (R01GM079383 and P30DA035778).

446 Author Contributions

- 447 J.W designed the experiment; BJ conducted the experiment and analyzed the data; all authors
- 448 discussed and wrote the paper.

449 Acknowledgements

- 450 The authors also thank the computing resources provided by the Center for Research Computing
- 451 (CRC) at University of Pittsburgh.
- 452

454 **References**

- 455 1. Ferreira LG, Dos Santos RN, Oliva G, Andricopulo AD (2015) Molecular docking and
 456 structure-based drug design strategies, Mol;20:13384-13421.
- 457 2. Huang SY, Zou X (2010) Advances and challenges in protein-ligand docking, Int J Mol
 458 Sci.;11:3016-3034.
- 459 3. Willett P, Barnard JM, Downs GM (1998) Chemical Similarity Searching, J. Chem. Inf.
 460 Comput Sci;38:983-996.
- 461 4. Mark A. Johnson GMM. Concepts and applications of molecular similarity. In: Sons J. W.
 462 (ed). New York, 1990, 393.
- 5. Stumpfe D, Hu Y, Dimova D, Bajorath J (2014) Recent Progress in Understanding Activity
 Cliffs and Their Utility in Medicinal Chemistry, J Med Chem;57:18-28.
- 465 6. He X, Man VH, Ji B, Xie XQ, Wang J (2019) Calculate protein-ligand binding affinities with
 466 the extended linear interaction energy method: application on the Cathepsin S set in the D3R
 467 Grand Challenge 3, J Comput Aided Mol Des;33:105-117.
- 468 7. Wang E, Sun H, Wang J et al. (2019) End-Point Binding Free Energy Calculation with
 469 MM/PBSA and MM/GBSA: Strategies and Applications in Drug Design, Chem
 470 Rev;119:9478-9508.
- 471 8. He X, Liu S, Lee T-S et al. (2020) Fast, Accurate, and Reliable Protocols for Routine
 472 Calculations of Protein–Ligand Binding Affinities in Drug Design Projects Using AMBER
 473 GPU-TI with ff14SB/GAFF, ACS Omega;5:4611-4619.
- 474 9. Wang L, Wu Y, Deng Y et al. (2015) Accurate and Reliable Prediction of Relative Ligand
- 475 Binding Potency in Prospective Drug Discovery by Way of a Modern Free-Energy Calculation
- 476 Protocol and Force Field, J Am Chem Soc;137:2695-2703.

- 477 10. Gaulton A, Hersey A, Nowotka M et al. (2016) The ChEMBL database in 2017, Nucleic Acids
 478 Res;45:D945-D954.
- 479 11. Liu J, Wang R (2015) Classification of current scoring functions, J Chem Inf Model;55:475480 482.
- 481 12. Berman HM, Westbrook J, Feng Z et al. (2000) The Protein Data Bank, Nucleic Acids
 482 Res;28:235-242.
- 483 13. Davies M, Nowotka M, Papadatos G et al. (2015) ChEMBL web services: streamlining access
 484 to drug discovery data and utilities, Nucleic Acids Res;43:W612-W620.
- 485 14. Gaulton A, Hersey A, Nowotka M et al. (2017) The ChEMBL database in 2017, Nucleic Acids
 486 Res;45:D945-d954.
- 487 15. Friesner RA, Banks JL, Murphy RB et al. (2004) Glide: A New Approach for Rapid, Accurate
 488 Docking and Scoring. 1. Method and Assessment of Docking Accuracy, J Med Chem;47:1739489 1749.
- 490 16. Wang Z, Sun H, Yao X et al. (2016) Comprehensive evaluation of ten docking programs on a
 491 diverse set of protein-ligand complexes: the prediction accuracy of sampling power and
- 492 scoring power, Phys Chem Chem Phys;18:12964-12975.
- 493 17. O'Boyle NM, Banck M, James CA, Morley C, Vandermeersch T, Hutchison GR (2011) Open
 494 Babel: An open chemical toolbox, J Cheminf;3:33.
- 495 18. Luccarelli J, Michel J, Tirado-Rives J, Jorgensen WL (2010) Effects of Water Placement on
 496 Predictions of Binding Affinities for p38α MAP Kinase Inhibitors, J Chem Theory
 497 Comput;6:3850-3856.
- 498 19. Michel J, Verdonk ML, Essex JW (2006) Protein-Ligand Binding Affinity Predictions by
- 499 Implicit Solvent Simulations: A Tool for Lead Optimization?, J Med Chem;49:7427-7439.

- 20. Pearlman DA, Charifson PS (2001) Are free energy calculations useful in practice? A
 comparison with rapid scoring functions for the p38 MAP kinase protein system, J Med
 Chem:44:3417-3423.
- 503 21. Jain AN, Nicholls A (2008) Recommendations for evaluation of computational methods, J
 504 Comput Aided Mol Des;22:133-139.
- 505 22. Sanner MF (1999) Python: a programming language for software integration and development,
 506 J Mol Graph Model;17:57-61.
- 507 23. O'Boyle NM, Banck M, James CA, Morley C, Vandermeersch T, Hutchison GR (2011) Open
 508 Babel: An open chemical toolbox, J Cheminformatics;3:33.
- 509 24. Huang N, Shoichet BK, Irwin JJ (2006) Benchmarking Sets for Molecular Docking, J Med
 510 Chem;49:6789-6801.
- 511 25. Mysinger MM, Carchia M, Irwin JJ, Shoichet BK (2012) Directory of Useful Decoys,
 512 Enhanced (DUD-E): Better Ligands and Decoys for Better Benchmarking, J Med
 513 Chem;55:6582-6594.
- 514 26. Trott O, Olson AJ (2010) AutoDock Vina: improving the speed and accuracy of docking with
 515 a new scoring function, efficient optimization, and multithreading, J Comput Chem;31:455516 461.
- 517 27. Irwin JJ, Shoichet BK (2005) ZINC--a free database of commercially available compounds for
 518 virtual screening, J Chem Inf Model;45:177-182.
- 519 28. Poli G, Galati S, Martinelli A, Supuran CT, Tuccinardi T (2020) Development of a
- 520 cheminformatics platform for selectivity analyses of carbonic anhydrase inhibitors, J Enzyme
- 521 Inhib Med Chem;35:365-371.

522	29. Shen H, Zamboni N, Heinonen M, Rousu J (2013) Metabolite Identification through Machine
523	Learning- Tackling CASMI Challenge Using FingerID, Metab;3:484-505.

- 524 30. Wang L, Wu Y, Deng Y et al. (2015) Accurate and Reliable Prediction of Relative Ligand
- 525 Binding Potency in Prospective Drug Discovery by Way of a Modern Free-Energy Calculation
- 526 Protocol and Force Field, J Am Chem Soc;137:2695-2703.
- 527 31. Lee T-S, Cerutti DS, Mermelstein D et al. (2018) GPU-Accelerated Molecular Dynamics and
- 528 Free Energy Methods in Amber18: Performance Enhancements and New Features, J Chem Inf
 529 Model;58:2043-2050.
- 530 32. Lee T-S, Hu Y, Sherborne B, Guo Z, York DM (2017) Toward Fast and Accurate Binding
- Affinity Prediction with pmemdGTI: An Efficient Implementation of GPU-Accelerated
 Thermodynamic Integration, J Chem Theory Comput;13:3077-3084.
- 533

535 FIGURES





Fig. 2 The comparison of RMSE, MAE, R² and PI values before and after the calibration for 11 drug receptors under the conditions of different fingerprint type and chemical similarity effect function. "*orig*" refers to original docking scores before the calibration.



546

Fig. 3 The changes of RMSE, MAE, R^2 and PI values between the calibrated docking scores and the original docking scores using the Glide docking scoring function under the conditions of different fingerprint type and chemical similarity effect function, S^p , where p takes a value of 1, 2, 3 or 4.





Fig. 4 The mean ROC curves of screening results before and after calibration of Glide docking



- *"cali"* and *"orig"* represent the calibrated and original docking scores, respectively.

558 TABLES

Receptor	Total Compounds	Activities (K _i /K _d)	Reference set	Validation set
CFX	7084	4521	581	144
D2R	9897	10473	635	161
ERK2	19729	1767	690	171
ER	6561	1452	191	51
MOR	7939	5796	362	110
VEGFR2	12497	1521	509	121
5HT2AR	6524	5130	619	150
A2AR	9015	7172	723	166
CB1	8907	5544	529	117
M1R	4858	2196	469	108
rRNA	79	77	57	14

Table 1. The information data of compounds for 11 receptors.

- **Table 2.** Mean RMSE (kcal/mol), MAE (kcal/mol), R² and PI before and after calibration using
- 575 the best hybrid scoring function (FP2 fingerprint with CSE=S⁴) for 11 receptors. "*cali*" and
- *"orig"* represent the calibrated and original docking scores, respectively.

	Receptor	Glide						AutoDock Vina									
	-	MAE RMSE		ISE	R ²		P	I	MAE		RMSE		R ²		PI		
		cali	orig	cali	orig	cali	orig	cali	orig	cali	orig	cali	orig	cali	orig	cali	orig
	CFX	1.50	2.10	1.89	2.60	0.52	0.18	0.72	0.41	1.40	1.96	1.76	2.35	0.50	0.12	0.72	0.33
	D2R EDV2	1.60	2.49	2.03	3.02	0.17	0.01	0.41	0.09	1.35	1.85	1.71	2.23	0.26	0.03	0.52	0.14
	ERK2 ED	0.88	1.20	1.16	1.55	0.57	0.35	0.81	0.74	0.6/	1.48	0.89	1./8	0.70	0.11	0.8/	0.54
	LN	1.20	2.01	2 33	2.45	0.09	0.30	0.82	0.05	1.30	1.70	1.70	2.17	0.31	0.14	0.71	0.30
	VEGFR2	1.71	1.80	2.08	2.33	0.39	0.02	0.64	0.51	1.02	1.43	1.32	1.80	0.58	0.42	0.83	0.66
	5HT2AR	1.28	2.38	1.61	2.88	0.39	0.05	0.63	0.23	1.27	1.59	1.62	1.98	0.32	0.11	0.58	0.36
	A2AR	1.32	1.88	1.74	2.32	0.46	0.08	0.68	0.29	1.30	1.89	1.70	2.33	0.44	0.03	0.67	0.14
	CB1	1.23	1.76	1.59	2.15	0.39	0.04	0.64	0.20	1.12	1.99	1.41	2.38	0.49	0.15	0.71	0.39
	M1R	1.36	1.83	1.72	2.27	0.49	0.14	0.72	0.33	1.18	1.62	1.48	1.99	0.58	0.31	0.77	0.54
	rRNA	0.70	1.86	0.85	2.35	0.58	0.08	0.74	0.19	0.63	1.88	0.80	2.10	0.70	0.26	0.78	0.48
	Average	1.34	2.05	1.69	2.53	0.44	0.14	0.66	0.34	1.15	1.74	1.47	2.11	0.50	0.16	0.71	0.39
578 579 580 581 582 583 584 585 586 587 588 589																	
590																	

Table 3. Mean AUC, $EF_{10\%}$ and $EF_{40\%}$ of screening results before and after calibration of Glide docking scores using the best hybrid scoring function (FP2 fingerprint with CSE=S⁴) for 11 receptors. "*cali*" and "*orig*" represent the calibrated and original docking scores, respectively. The average numbers of compounds allocated in the active and inactive sets are shown in the table. 95% CI for each metrics is displayed in the parenthesis.

Receptor	EF	10%	EF	40%	A	UC	Actives	Inactives
	cali	orig	cali	orig	cali	orig		
CFX	2.24 (0.13)	1.85 (0.15)	1.72 (0.08)	1.46 (0.08)	0.85 (0.02)	0.72 (0.03)	84 (6)	62 (3)
D2R	1.96 (0.31)	1.34 (0.15)	1.48 (0.14)	1.14 (0.12)	0.68 (0.05)	0.53 (0.05)	55 (4)	104 (6)
ERK2	7.73 (0.43)	6.26 (0.41)	2.44 (0.07)	2.40 (0.09)	0.97 (0.02)	0.93 (0.02)	20 (2)	151 (7)
ER	2.67 (0.25)	1.91 (0.42)	2.24 (0.09)	1.90 (0.16)	0.95 (0.01)	0.86 (0.03)	18 (3)	32 (4)
MOR	1.52 (0.33)	1.19 (0.31)	1.36 (0.13)	1.13 (0.10)	0.72 (0.05)	0.59 (0.04)	45 (4)	51 (8)
VEGFR2	2.15 (0.12)	2.18 (0.10)	1.75 (0.06)	1.64 (0.06)	0.86 (0.03)	0.79 (0.03)	59 (6)	70 (4)
5HT2AR	1.81 (0.15)	1.20 (0.18)	1.54 (0.08)	1.13 (0.04)	0.81 (0.02)	0.58 (0.02)	77 (5)	85 (5)
A2AR	2.45 (0.13)	1.66 (0.22)	1.87 (0.06)	1.37 (0.05)	0.86 (0.01)	0.66 (0.02)	67 (5)	109 (6)
CB1	2.09 (0.15)	1.67 (0.21)	1.83 (0.08)	1.32 (0.12)	0.84 (0.02)	0.63 (0.03)	54 (4)	76 (4)
M1R	2.39 (0.24)	2.33 (0.28)	1.90 (0.08)	1.33 (0.11)	0.87 (0.02)	0.67 (0.03)	42 (5)	73 (6)
rRNA	11.50 (2.99)	0.00	2.59 (0.20)	2.59 (0.20)	1.00	0.86 (0.03)	1	13 (2)

596

597

Table 4. Mean RMSE (kcal/mol), MAE (kcal/mol), R² and PI before and after calibration of

600 Glide docking scores for 11 receptors by giving similarity a range with lower and upper bound in

601 diverse library. "Cali" and "Orig" represent the calibrated and original docking scores,

602 respectively. 95% CI for each metrics is displayed in the parenthesis.

Receptor	Tc range	MA	4E	RM	ISE	F	R ²	PI		
		Cali	Orig	Cali	Orig	Cali	Orig	Cali	Orig	
CFX	[0.30,0.45]	1.96	2.01	2.40	2.52	0.20	0.20	0.46	0.45	
		(0.06)	(0.07)	(0.07)	(0.07)	(0.04)	(0.04)	(0.04)	(0.04)	
D2R	[0.30,0.40]	1.92	2.47	2.41	3.01	0.01	0.01	0.10	0.09	
		(0.05)	(0.05)	(0.07)	(0.06)	(0.01)	(0.01)	(0.04)	(0.04)	
ERK2	[0.30,0.40]	1.09	1.19	1.38	1.54	0.43	0.36	0.77	0.76	
		(0.03)	(0.04)	(0.07)	(0.05)	(0.05)	(0.04)	(0.04)	(0.04)	
ER	[0.20,0.35]	1.79	2.11	2.15	2.53	0.43	0.37	0.69	0.65	
		(0.09)	(0.08)	(0.07)	(0.06)	(0.06)	(0.05)	(0.05)	(0.04)	
MOR	[0.35,0.40]	2.19	3.11	2.70	3.77	0.07	0.03	0.23	0.14	
		(0.10)	(0.15)	(0.10)	(0.17)	(0.02)	(0.02)	(0.06)	(0.07)	
VEGFR2	[0.25,0.40]	1.89	1.78	2.34	2.33	0.25	0.24	0.51	0.49	
		(0.08)	(0.08)	(0.11)	(0.09)	(0.05)	(0.04)	(0.05)	(0.05)	
5HT2AR	[0.30,0.45]	1.69	2.36	2.09	2.87	0.09	0.07	0.25	0.24	
		(0.07)	(0.09)	(0.08)	(0.09)	(0.04)	(0.03)	(0.06)	(0.08)	
A2AR	[0.35,0.40]	1.84	1.84	2.29	2.27	0.13	0.09	0.36	0.29	
		(0.08)	(0.09)	(0.08)	(0.09)	(0.03)	(0.03)	(0.05)	(0.06)	
CB1	[0.25, 0.40]	1.65	1.78	2.08	2.19	0.09	0.04	0.30	0.21	
		(0.05)	(0.05)	(0.05)	(0.06)	(0.03)	(0.02)	(0.06)	(0.06)	
M1R	[0.35,0.40]	1.80	1.83	2.20	2.25	0.24	0.13	0.47	0.32	
		(0.10)	(0.07)	(0.12)	(0.08)	(0.04)	(0.03)	(0.04)	(0.04)	
rRNA	[0.30,0.35]	1.64	1.78	1.93	2.24	0.14	0.11	0.41	0.22	
		(0.11)	(0.23)	(0.11)	(0.25)	(0.06)	(0.06)	(0.10)	(0.18)	
Average		1.77	2.02	2.18	2.50	0.19	0.15	0.41	0.35	

Table of Contents

