

Selection of In Silico Drug Screening Results by Using Universal Active Probes (UAPs)

Yoshifumi Fukunishi,^{*,†,‡} Kazuki Ohno,^{§,||} Masaya Orita,^{§,||} and Haruki Nakamura^{†,⊥}

Biomedical Information Research Center (BIRC), National Institute of Advanced Industrial Science and Technology (AIST), 2-3-26, Aomi, Koto-ku, Tokyo 135-0064, Japan, Pharmaceutical Innovation Value Chain, BioGrid Center Kansai, 1-4-2 Shinsenri-Higashimachi, Toyonaka, Osaka 560-0082, Japan, Japan Biological Informatics Consortium (JBIC), 2-41-6, Aomi, Koto-ku, Tokyo 135-0064, Japan, Chemistry Research Laboratories, Drug Discovery Research, Astellas Pharma Incorporated, 21 Miyukigaoka, Tsukuba, Ibaraki 305-8585, Japan, and Institute for Protein Research, Osaka University, 3-2 Yamadaoka, Suita, Osaka 565-0871, Japan

Received March 17, 2010

We developed a new method that uses a set of drug-like compounds to select reliable in silico drug screening results. If some active compounds are known, the screening results that rank these active compounds at the top should be reliable. If no active compound is known, how to select the result is in question. We propose a concept of a set of “universal active probes” (UAPs), which is a set of small active compounds that bind to different kinds of proteins. We found that the hit ratio of the true active compounds in in silico screening shows positive correlation to that of the UAPs, probably because UAPs form a set of drug-like compounds. Thus, if the UAPs were added to the compound library, the screening result that shows a high hit ratio of the UAPs could give reliable actual hit compounds for the target protein. We examined this method for several targets and found this idea useful.

1. INTRODUCTION

Structure-based in silico drug screening methods are composed of several techniques, including the modeling of a target protein, the generation of a chemical compound library, a protein–compound docking program, and the postprocessing of the docking poses and docking scores obtained by protein–compound docking. Even if the atomic coordinates of a target protein are experimentally determined, their multiple structures are often useful for a protein–compound docking study to consider the induced fitting due to the compound binding.^{1–3} The protein–compound docking program can generate a drug screening result for each protein structure so that many screening results are obtained from the many protein structures that a protein simulation program generates. The problem is how to select the best of the many screening results.^{1–3}

In silico drug screening results by structure-based screening methods depend heavily on the target protein's three-dimensional (3D) structure.^{4,5} Even if target protein structures determined by X-ray crystallography experiments are used, in silico screening succeeds in providing good database enrichment in approximately half of the cases and fails in hit compound prediction in the other half.⁴ In some cases, the prediction results are much worse than the results obtained by random screening. We previously reported that the model structures obtained by the MD simulations would provide poorer enrichment than the crystal structures, while

a holo structure could be obtained from an apo structure by the MD simulation.³ This problem could not be overcome by improving the docking program. For a target protein, in silico screening results depend heavily on which protein–compound docking program is used.⁴ There is at the moment no perfect docking program, and the screening results depend on the combination of the docking program used and the target protein structure. In particular, even a slight structural change around the binding site will sometimes have a large effect on the docking scores.^{6–10}

Ensemble docking is becoming the new trend of the structure-base in silico drug screening recently.^{11–13} In the ensemble docking, multiple target protein structures are prepared to consider the flexibility of the protein structure, and the compounds of the library are docked to these multiple structures. This approach is effective in some cases, but the problem is how to select the target structure that will give a good screening result among the multiple protein structures. If some active compounds were known for the target protein, the structure that shows the good hit ratio of these known active compounds should be selected as the suitable protein structure for the docking screening study. If no active compound was known for the target protein, it is difficult to select the suitable protein structure for the docking screening study.

In the present study, we showed how to select the target structure that will give the good screening result among the multiple protein structures without knowledge of known active compounds. Several kinds of compounds were prepared for the target proteins, which were soluble. These compounds are the ligands of the target protein, the G-protein coupled receptor (GPCR), and the other soluble proteins. Interestingly, we found that the target proteins can also bind the ligands of the other proteins as well as their own ligands.

* Corresponding author. Telephone: +81-3-3599-8290. Fax +81-3-3599-8099. E-mail: y-fukunishi@aist.go.jp.

[†] National Institute of Advanced Industrial Science and Technology.

[‡] Pharmaceutical Innovation Value Chain, BioGrid Center Kansai.

[§] Japan Biological Informatics Consortium.

^{||} Astellas Pharma Incorporated.

[⊥] Osaka University.

In other words, drug-like compounds show a higher affinity to many kinds of proteins than nondrug-like compounds do. We called such drug-like compounds “universal active probes” (UAPs). We could assess the quality of the *in silico* screening using the hit ratio of the UAPs even without knowing the active compounds of the target protein.

2. METHODS

Multiple 3D structures of the target protein were generated by molecular dynamics (MD) simulation, and the structure-based drug screening method provided a screening result for each 3D structure of the target protein. The structure-based *in silico* drug screening was performed by the multiple target screening (MTS) method using the SievGene protein–compound docking program (see Appendix A).^{14–16} Then, the screening results (database enrichment) for true active compounds and those of the UAPs were compared to each other. Each screening result was evaluated by the area under the database enrichment curve (AUC). The AUC values are calculated by

$$\text{AUC} = \int_{0\%}^{100\%} f(x)dx \quad (1)$$

where x and $f(x)$ are the percentages of compounds that are selected from the total compound library and the database enrichment curve, respectively. A higher AUC value corresponds to better database enrichment, and the AUC value is always more than zero and less than 100. For the random screening, $\text{AUC} = 50$. The AUC value for the true active compounds is the frequency of finding the true active compounds. For the AUC of the UAP, the database enrichment curve represents the finding frequency of the UAPs instead of the true active compounds.

2.1. Universal Active Probe. A UAP is a small active compound of an arbitrary target protein. We prepared three types of UAPs: (i) the ligands of GPCRs, (ii) those of many soluble proteins, and (iii) those of the directory of useful decoy (DUD).¹⁷

The GPCR ligands are exactly the same as the compounds used in our previous study.¹⁸ This UAP set, which we call UAP_GPCR, consisted of 73 compounds, including 10 antagonists of histamine H1 receptor, 12 agonists and 13 antagonists of adrenaline β -receptor, 8 agonists and 9 antagonists of serotonin receptor, and 6 agonists and 15 antagonists of dopamine D2 receptor.

The second UAP set consisted of the 175 ligands of the protein–complex structure listed in Appendix B, except for the ligands of the target proteins. We call this set UAP_PDB.

Our third set of UAPs was selected from the active compounds of the DUD release 2 (<http://dud.docking.org/r2/>). Details about the DUD decoy set were given in an earlier paper.¹⁷ There were 3961 active compounds of the 39 target proteins. Active compounds with mass weight (MW) ≥ 250 Da and MW ≤ 350 Da were selected. Clustering analysis was performed for these DUD active compounds based on the Pipeline Pilot FCFP_4 molecular descriptor,^{19,20} and the analysis generated 148 clusters. A representative compound of each cluster was selected as the UAP, yielding a total of 148 UAPs. We call this set the UAP_DUD.

The average numbers of atoms were 44.80, 49.07, and 38.39 for UAP_GPCR, UAP_PDB, and UAP_DUD, respec-

tively. The average numbers of rings were 2.697, 2.160, and 2.993 atoms for the UAP_GPCR, UAP_PDB, and UAP_DUD, respectively. Thus the average numbers of atoms and rings were similar among these UAPs.

The 3D coordinates of the UAPs were generated by myPresto (myPresto; <http://medals.jp/myPresto/license.html> and myPresto Management; http://presto.protein.osaka-u.ac.jp/myPresto4/index_e.html). We used the general AMBER force field (GAFF),²¹ and the molecular topology files were generated by tpgeneL/myPresto. The energy of the coordinates of small molecules was optimized by Cosgene/myPresto.²² The atomic charges were calculated by the Gasteiger method of Hgene/myPresto.^{23,24}

2.2. Protein Sets. Our screening methods are based on a protein–compound affinity matrix. We, therefore, must prepare a set of proteins including target proteins. An individual protein set was prepared for each target structure. Each protein set consisted of a basic protein set and the target structure itself. All structures of the basic protein set were crystal structures.

To evaluate our method, we performed a protein–compound docking simulation based on the soluble protein structures registered in the Protein Data Bank (PDB). The protein sets (basic protein sets) used are exactly the same as those used in our previous study.^{25,26} Here we briefly describe the data set again.

The protein–ligand complex structures were suitable for the docking study, since the ligand pockets were clearly determined. A total of 180 proteins were selected from the PDB, 142 complexes were selected from the database used in the evaluation of GOLD and FlexX,²⁷ and the other 38 complexes were selected from the PDB. The 142 protein data set contains a rich variety of proteins and compounds whose structures have all been determined by high-quality experiments with a resolution of less than 2.5 Å.²⁷ Almost all the atomic coordinates are supplied except for those of the hydrogen atoms, and the atomic structures around the ligand pockets are reliable. Thus, this data set was used in the clustering analysis of proteins and *in silico* screening. From the original data set, the complexes containing a covalent bond between the protein and ligand were removed, since our docking program cannot perform protein–ligand docking when a covalent bond exists between the protein and the ligand. The other 38 structures include the human immunodeficiency virus protease-1, cyclooxygenase-2 (COX2), and glutathione S-transferase. The PDB identifiers are summarized in Appendix B. All water molecules and cofactors were removed from the proteins, and all missing hydrogen atoms were added to form the all-atom models of the proteins.

2.3. Preparation of Target Protein Structures and Their Inhibitors. Cyclooxygenase-2 (COX2), AmpC β -lactamase (AMPC), factor Xa (FXA), and thrombin (THR) were selected for the test to validate the current method. Six target protein structures (PDB IDs: 1cx2, 1pxx, 3pgh, 4cox, 5cox, and 6cox) were selected for COX2. Two target protein structures were selected for each of the other proteins: 2pu2 and 2r9x for AMPC, 2w26 and 3ens for FXA, and 2pks and 2zgp for THR. The numbers of prepared inhibitors for COX2, AMPC, FXA, and THR were 9, 10, 10, and 12, respectively. These inhibitors were exactly the same as those used in our previous study, and they are summarized in the Supporting

Information.²⁵ The 3D data were prepared by myPresto in the same manner as the UAP data.

The model structures generated from the PDB structure given by MD in explicit water were prepared as follows. All target proteins were prepared without ligands (apo structure). The force fields and charges of protein atoms were originated from AMBER parm99.²⁸ The whole structure of each protein was embedded in a sphere of TIP3P²⁹ water (CAP water), including ion particles of 0.1% Na⁺ and Cl⁻, in order to neutralize the total charge of the systems. The center of the sphere was set at the mass center of the protein. The radii were 52.502, 39.016, 39.016, 43.609, 44.563, 43.592, 39.741, 39.741, 62.076, 51.578, 51.978, 52.341, 52.0958 Å for 1pxx, 2pks, 2pu2, 2r9x, 2w26, 2zgp, 3ens, 3pgh, 4cox, 5cox, and 6cox, respectively. Before MD calculations were performed for the entire system, a MD calculation for only the solvent parts (solvent water and counterions) was performed with the protein, ligand, and metal ion coordinates fixed, so as to bring the solvent parts sufficiently close to an equilibrium state. The SHAKE method was used to constrain covalent bonds between heavy and hydrogen atoms in any molecule in the system.³⁰ MD simulations of the entire system were performed using 1.5 fsec time steps with the temperature set at 310 K and the fast multipole method³¹ being used to calculate the Coulombic interaction. The cutoff distance of the van der Waals interaction was 10.0 Å. The MD simulations were performed by using cosgene/myPresto.²² After equilibration steps of 500 psec, the protein coordinates were sampled every 100 psec. Finally, we obtained 12 structures for each target protein, including 10 sampled structures, 1 energy-minimized structure, and the initial structure of the target protein. The average root-mean-square deviation (rmsd) value of heavy atoms of the 12 structures was calculated for these 11 target proteins. The average rmsd values were 2.241, 2.145, 1.845, 1.765, 2.100, 2.074, 3.768, 2.385, 2.497, 2.570, 2.479 Å for 1pxx, 2pks, 2pu2, 2r9x, 2w26, 2zgp, 3ens, 3pgh, 4cox, 5cox, and 6cox, respectively. These 12 structures were used as the structural ensemble of each target protein in the following analysis.

2.4. Preparation of Chemical Compound Libraries. We prepared three decoy compound libraries. They were the random library provided by the Coelacanth Corporation (Decoy 1), the randomly selected compounds from the LigandBox database³² (Decoy 2), and the DUD (Decoy 3).

One decoy set (Decoy 1) was the Coelacanth chemical compound library (Coelacanth Corporation, East Windsor, NJ), which is a random library consisting of 11 050 potential negative compounds. The Coelacanth decoy set was used for all targets. The second decoy set (Decoy 2) was the randomly selected 10⁴ compounds of the LigandBox database. The third decoy set (Decoy 3) was the decoy set of the DUD for each target protein.¹⁷ Specific DUD decoy sets were prepared for each target. The numbers of compounds in the DUD decoy sets for COX2, AMPC, FXA, and THR were 13 289, 786, 5745, and 2456, respectively. The average number of atoms, the average number of heavy atoms, and the average mass weights of these decoy sets are summarized in Table 1. Usually only one hit compound was found out of 10⁴ randomly selected compounds; thus, we expected that there were few, if any, hit compounds among these 10⁴ compounds. The 3D coordinates of the 11 050 chemical

Table 1. Average Number of Atoms and Heavy Atoms and Average Mass Weight (Da) of Decoy Set and True Active Compounds for Each Target Protein

		COX2	AMPC	FXA	THR
Coelacanth decoy (Decoy 1)	no. of atoms	63.6	63.6	63.6	63.6
	no. of heavy atoms	30.9	30.9	30.9	30.9
	mass weight	423.0	423.0	423.0	423.0
LigandBox decoy (Decoy 2)	no. of atoms	37.2	37.2	37.2	37.2
	no. of heavy atoms	20.4	20.4	20.4	20.4
	mass weight	289.2	289.2	289.2	289.2
DUD decoy (Decoy 3)	no. of atoms	40.2	30.2	54.5	57.6
	no. of heavy atoms	25.3	20.8	32.6	32.3
	mass weight	364.1	304.8	455.7	453.2
Original ligand	no. of atoms	35.0	30.0	55.6	66.1
	no. of heavy atoms	22.0	20.5	33.4	34.7
	mass weight	316.9	315.6	465.9	495.4

compounds of the Coelacanth chemical compound library (Decoy 1) were generated by the Concord program (Tripos, St. Louis, MO) from the 2D Sybyl SD files provided by the Coelacanth Chemical Corporation. The 3D coordinates of the 10 000 compounds of the LigandBox library (Decoy 2) were prepared by using the programs of myPresto with the GAFF, and the detail was described in the earlier publication.³² The 3D coordinates of the DUD (Decoy 3) were retrieved from the DUD Web site (<http://dud.docking.org/>).

3. RESULTS

Figure 1 shows the examination scheme. The compound library consists of the true active compounds, the UAPs, and the decoy compounds. For the target protein, 12 structures were prepared by using the MD simulation. Structure-based in silico drug screening was performed using the MTS method. The protein–compound docking procedure was exactly the same as that reported in our previous work.¹³ For flexible docking, the SievGene program generated up to 100 conformers for each compound. We obtained a compound list for each target protein structure. The database enrichment curve was calculated for the true active compounds, and the area under the curve (AUC_{true}) was calculated

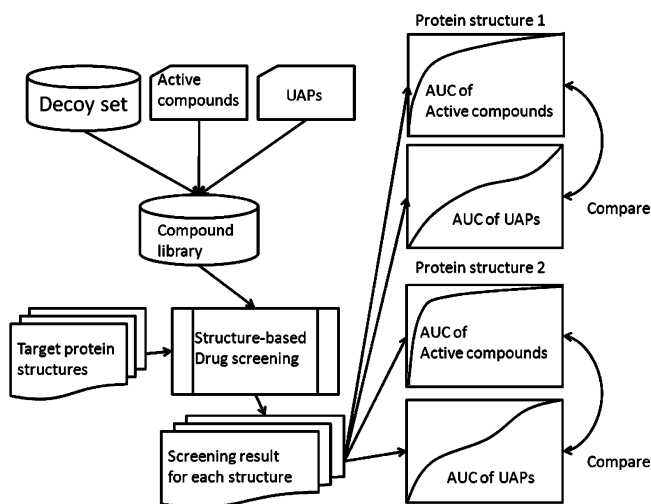


Figure 1. Schematic representation of the procedure used in the current study.

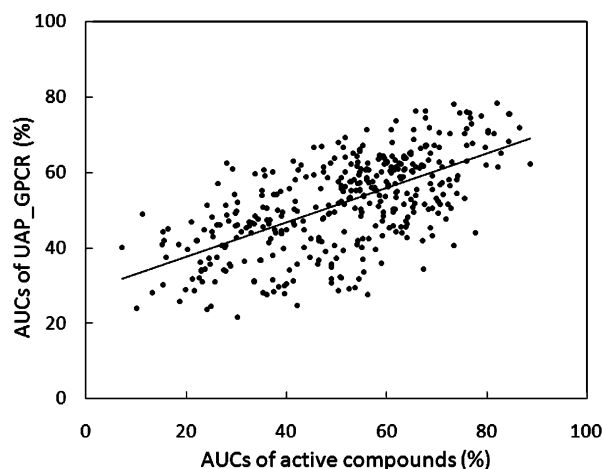


Figure 2. AUC_{true} vs AUC_{GPCR} . The least-squares fitting showed that $AUC_{true} = 2.17 * AUC_{GPCR} + 61.71$ and the correlation coefficient = 0.63. The solid line represents the result by the least-squares fitting.

Table 2. Correlation Coefficients (R) between the AUC_{true} and AUC_{UAP} Values

UAP	correlation coefficient (R) ^a	R		
		Decoy1 ^b	Decoy2 ^b	Decoy3 ^b
UAP_GPCR	0.6295	0.6981	0.5295	0.5834
UAP_DUD	0.5151	0.6173	0.3622	0.4477
UAP_PDB	0.4165	0.7468	0.6161	0.6347
UAP_min	0.6927	0.7582	0.6190	0.6386
UAP_avg	0.6608	0.7297	0.5723	0.6072
UAP_max	0.6040	0.6753	0.5001	0.5529

^a The R value between AUC_{true} and the AUC_{UAP} of a total of 396 screening results (11 targets \times 12 structures \times 3 decoy sets).

^b The R value between AUC_{true} and the AUC_{UAP} of a total of 132 screening results (11 targets \times 12 structures).

as a measure of the screening result. Also, the database enrichment curve was calculated for the UAPs, and the AUC (AUC_{UAP}) was also calculated. Finally, the AUC_{true} and the AUC_{UAP} values were compared, and the correlation coefficient between them was calculated. We also provided the Supporting Information to explain the scheme.

Figure 2 shows the correlations between the AUC_{true} and the AUC_{UAP} of UAP_GPCR. There are 396 data (11 target proteins \times 12 protein structures \times 3 decoy sets) plotted. The correlation coefficients are summarized in Table 2. The AUC_{true} value shows positive correlation to the AUC_{UAP} value, but the correlation is somewhat weak. The correlation between the AUC_{true} and the AUC_{GPCR} was stronger than the other correlations.

The average AUC_{true} values of the top 10, 20, and 30 results predicted by the UAP_GPCR were 76.43, 73.41, and 72.47, respectively. The average AUC_{true} values of the worst 10, 20, and 30 results predicted by the UAP_GPCR were 31.79, 33.01, and 34.90, respectively. Thus, the UAP_GPCR should be a good guide to select a good screening result.

Table 3 shows the relationships between AUC_{true} and AUC_{GPCR} for each target protein and decoy set. The correlation coefficients were calculated for the 12 protein structures of each target protein. In many cases, the correlations were positive. For some target proteins, the correlations were negative. This means that the result was not desirable; the structure with the high AUC_{GPCR} gave the low

AUC_{true} . The individual correlation coefficients were smaller than the overall correlation coefficient.

We also examined the consensus of the UAPs. For a single screening result, three AUCs were obtained for UAP_GPCR, UAP_DUD, and UAP_PDB. The average value of these three AUCs was defined as AUC_{avg} . The minimum and maximum values of the three AUCs were defined as AUC_{min} and AUC_{max} , respectively. These results are summarized in Table 2, and Figure 3 shows the correlations between the AUC_{true} and the AUC_{UAP} of UAP_min. The AUC_{min} gave the best correlation coefficient among these six R values. The AUC_{avg} gave the second-best correlation coefficient. These results showed that the consensus of the UAPs, especially AUC_{min} , was useful for predicting a good screening result.

Table 4 shows the relationships between the AUC_{true} and the AUC_{min} for each target protein and decoy set. The correlation coefficients were calculated for the 12 protein structures of each target protein. In many cases, the correlations were positive, though they were negative for some target proteins. The individual correlation coefficients were smaller than the overall correlation coefficient. These results were better than the data in Table 3.

The average AUC_{true} values of the top 10, 20, and 30 results predicted by the UAP_min were 76.67, 74.53, and 73.19, respectively. The average AUC_{true} values of the worst 10, 20, and 30 results predicted by the UAP_min were 31.32, 29.92, and 31.17, respectively. Thus, UAP_min should be useful for selecting a good screening result.

Table 5 shows the similarities between compounds. If the active compounds are more similar to the UAP than the compounds of the decoy set, then our analysis is not so useful. The similarities between compounds were calculated by the molecular dynamics maximum overlap (MD-MVO) method, which evaluates the volume overlap between two molecules considering the atomic charge.^{33,34} A MD-MVO score of -1 indicates a perfect match where the volume overlaps completely, and a score of 0 indicates a perfect mismatch. Table 5 shows that the active compounds were similar to each other and that the UAPs were also similar to each other. The average similarity between the active compounds and the UAPs was almost the same as that between the active compounds and the decoy compounds. These results showed that the UAPs were not similar to the active compounds. In other words, the results shown in Figures 2 and 3 and in Tables 2–4 are not trivial, and the UAPs could be applied to other many target proteins.

Figure 1 also shows how to select the protein structure that will give the high hit ratio among multiple target protein structures in the ensemble docking study by using the UAP. The procedure is as follows:

- Step 1: The compound library includes the UAPs.
- Step 2: Multiple target protein structures are prepared. These structures can be generated by the MD simulation and also are collected from the PDB.
- Step 3: The structure-based in silico drug screening is performed for these multiple target protein structures to give multiple screening results.
- Step 4: The screening result (and the target protein structure) that gives the high hit ratio of the UAP is selected as the final result (and the suitable structure for docking screening study).

Table 3. AUC_{true} and AUC_GPCR Values (%) for Each Target Protein

target protein	decoy set	AUC _{true} ^a	AUC_GPCR ^a	AUC _{true} with highest AUC_GPCR ^b	AUC _{true} with lowest AUC_GPCR ^c	correlation coefficient (R) ^d
1pxx	Decoy1	70.35	71.77	73.38	52.69	0.705
	Decoy2	56.76	61.83	47.08	59.77	-0.182
	Decoy3	55.96	62.05	64.92	51.77	0.356
2pks	Decoy1	37.40	46.90	38.92	37.25	-0.184
	Decoy2	32.79	45.29	63.42	26.17	0.409
	Decoy3	33.40	45.80	64.83	23.75	0.430
2pu2	Decoy1	53.25	40.48	71.90	51.00	0.682
	Decoy2	50.33	34.62	77.80	39.60	0.634
	Decoy3	42.56	28.84	73.40	30.30	0.832
2r9x	Decoy1	32.63	33.59	57.90	36.20	0.485
	Decoy2	42.49	41.74	41.60	54.40	-0.219
	Decoy3	27.53	29.14	41.30	10.20	0.276
2w26	Decoy1	55.51	56.98	55.30	56.50	-0.230
	Decoy2	46.39	47.98	68.70	22.90	0.707
	Decoy3	44.81	45.64	33.70	35.00	0.374
2zgp	Decoy1	37.17	50.03	29.25	49.00	-0.759
	Decoy2	31.51	46.19	50.92	20.25	0.886
	Decoy3	31.66	46.18	35.17	49.08	0.192
3ens	Decoy1	57.98	53.87	51.40	62.80	-0.153
	Decoy2	53.71	49.46	51.00	46.40	-0.283
	Decoy3	52.40	47.46	38.00	55.20	-0.164
3pgh	Decoy1	72.24	66.16	80.23	58.38	0.304
	Decoy2	61.34	59.97	67.92	62.00	0.098
	Decoy3	59.56	54.37	65.31	61.23	0.254
4cox	Decoy1	66.57	64.82	76.62	57.00	0.468
	Decoy2	62.09	54.66	56.31	68.69	-0.282
	Decoy3	57.95	56.59	61.31	52.62	-0.040
5cox	Decoy1	70.33	69.91	82.08	62.92	0.676
	Decoy2	62.15	61.83	51.67	72.67	-0.391
	Decoy3	59.21	62.06	73.00	58.00	0.774
6cox	Decoy1	62.37	64.43	86.46	57.54	0.707
	Decoy2	63.38	59.23	69.46	74.08	-0.088
	Decoy3	54.92	58.25	35.46	41.69	0.114
average		51.48	52.06	58.66	48.40	0.224

^a Average AUC value of the AUC values for 12 protein structures for a target protein. ^b AUC_{true} value of a screening result giving the highest AUC_GPCR value. ^c AUC_{true} value of a screening result giving the lowest AUC_GPCR value. ^d R value among the AUC_{true} values and the AUC_GPCR values for 12 protein structures for a target protein.

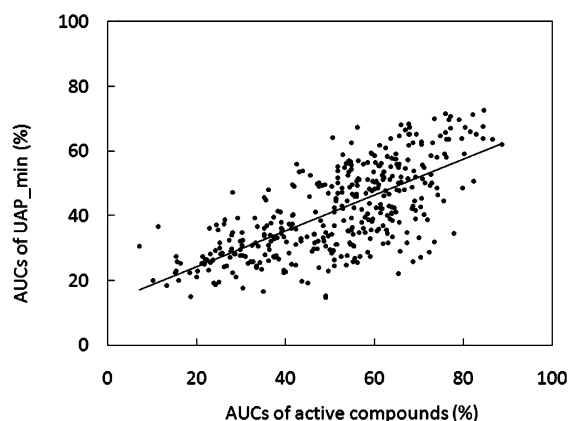


Figure 3. AUC_{true} vs AUC_{min}. The least-squares fitting showed that the AUC_{true} = 1.80 * AUC_{min} = 23.47 and the correlation coefficient = 0.69. The solid line represents the result by the least-squares fitting.

4. DISCUSSION

The *R* value between the AUC_{true} and the AUC_GPCR is bigger than that between the AUC_{true} and the AUC_DUD. We examined the difference between the UAP_GPCR and the UAP_DUD. The average number of atoms of UAP_GPCRs was 44.80 and that of the UAP_DUDs was 38.39. The average number of rings of UAP_GPCRs was 2.697 and that of the UAP_DUDs was 2.993. The difference

between these numbers was small. It is difficult to find the difference between the UAP_GPCR and the UAP_DUD. One of the differences is that the UAP_GPCRs show strong affinity to the target GPCR, and the UAP_DUDs are hit-level compounds that should show weak affinity to the target protein. In other words, the UAP_GPCRs are more drug-like than the UAP_DUDs. The drug-likeness of the compound could be important for selecting the UAP. Also, the *R* value between the AUC_{true} and the AUC_GPCR is bigger than that between the AUC_{true} and the AUC_PDB. The same as with the UAP_DUD, the difference between the UAP_GPCR and the UAP_PDB was not so clear. The ligands of the protein–compound complexes registered in the PDB were not necessarily drugs, and the UAP_GPCR could be more drug-like than the UAP_PDB.

The hit ratio (AUC) of the true active compounds showed correlation to that of the UAPs. This result suggests that there should be a common feature between the true active compound and the UAP. The recent study suggests that there could be a common structural feature among the ligand-binding pockets and that a clustering analysis of these pockets is possible.³⁵ One possible explanation of our result is that the active compounds of proteins have a common structural feature that reflects the common structural feature of ligand-

Table 4. AUC_{true} and AUC_{min} Values (%) for Each Target Protein

target protein	decoy set	AUC _{true} ^a	AUC _{min} ^a	AUC _{true} with highest AUC _{min} ^b	AUC _{true} with lowest AUC _{min} ^c	correlation coefficient (<i>R</i>) ^d
1pxx	Decoy1	70.35	63.01	73.38	52.69	0.84
	Decoy2	56.76	52.80	64.77	71.08	0.02
	Decoy3	55.96	49.46	59.46	54.92	0.33
2pks	Decoy1	37.40	34.15	35.25	30.33	0.12
	Decoy2	32.79	30.33	63.42	28.00	0.58
	Decoy3	33.40	31.08	64.83	23.75	0.63
2pu2	Decoy1	53.25	34.37	69.20	51.00	0.77
	Decoy2	50.33	28.72	77.80	28.90	0.88
	Decoy3	42.56	23.85	73.40	18.70	0.92
2r9x	Decoy1	32.63	29.34	57.90	23.10	0.79
	Decoy2	42.49	34.43	33.00	25.40	0.16
	Decoy3	27.53	24.00	50.30	13.20	0.77
2w26	Decoy1	55.51	39.71	41.50	65.80	-0.63
	Decoy2	46.39	30.37	70.30	43.60	0.22
	Decoy3	44.81	29.10	69.10	35.00	0.07
2zgp	Decoy1	37.17	34.01	28.08	49.00	-0.67
	Decoy2	31.51	30.94	61.67	18.50	0.80
	Decoy3	31.66	29.78	61.25	49.08	0.25
3ens	Decoy1	57.98	38.24	51.40	54.60	-0.07
	Decoy2	53.71	32.32	63.80	49.60	0.29
	Decoy3	52.40	31.19	64.30	45.00	0.32
3pgh	Decoy1	72.24	61.91	82.77	58.38	0.38
	Decoy2	61.34	54.14	67.92	69.08	0.25
	Decoy3	59.56	49.48	62.38	58.69	0.24
4cox	Decoy1	66.57	54.70	75.92	58.38	0.57
	Decoy2	62.09	51.16	67.46	68.69	-0.38
	Decoy3	57.95	45.29	53.54	57.15	-0.15
5cox	Decoy1	70.33	61.81	76.00	62.92	0.66
	Decoy2	62.15	51.74	52.92	70.50	-0.23
	Decoy3	59.21	53.04	73.00	51.00	0.56
6cox	Decoy1	62.37	56.26	65.85	60.69	0.68
	Decoy2	63.38	55.17	69.46	74.23	-0.10
	Decoy3	54.92	48.99	82.23	41.69	0.50
average		51.48	41.66	62.53	47.35	0.314

^a Average AUC value of the AUC values for 12 protein structures for a target protein. ^b AUC_{true} value of a screening result giving the highest AUC_{min} value. ^c AUC_{true} value of a screening result giving the lowest AUC_{min} value. ^d *R* values among the AUC_{true} values and the AUC_{min} values for 12 protein structures for a target protein.

Table 5. Average Similarity between Compounds^a

similarity	MVO _{score}
ACT_ACT	
AMPC_AMPC	-0.621
COX2_COX2	-0.529
FXA_FXA	-0.477
THR_THR	-0.565
UAP_UAP	
GPCR_GPCR	-0.565
PDB_PDB	-0.329
DUD_DUD	-0.553
ACT_UAP	
ACT_GPCR	-0.401
ACT_PDB	-0.375
ACT_DUD	-0.397
ACT_DECOY	
ACT_Decoy1	-0.375
ACT_Decoy2	-0.391
ACT_Decoy3	-0.363

^a ACT_ACT represents the average similarity among the true active compounds. UAP_UAP represents the average similarity among the UAPs. ACT_UAP represents the average similarity among the true active compounds and the UAPs. ACT_DECOY represents the average similarity among the true active compounds and the decoy compounds. "ACT" represents all the true active compounds.

binding pockets. Such common structural features of active compounds could be related to a sort of drug-likeness.³⁶

The overall AUC_{true} values of the current study were not quite good. There have been many reports about the relationship between target modeling methods and hit ratios.^{6–10} Some reports suggested that the hit ratios depend on the structural changes in the ligand-binding regions.^{9,10} McGovern et al.⁶ reported that the holo crystal structures would give better enrichments than the apo crystal structures and that the apo structures would give better enrichments than the homology modeled structures. Our previous study showed that the modeled structure by the MD simulation gave poorer enrichment than the holo crystal structure in many cases.³ Because our MD simulation generated only the apo structures, the current result is not good but still reasonable.

Table 1 shows that the compounds of the Coelacanth decoy set are larger than the true active compounds of the target, the DUD decoy set is second largest, and the compounds of the LigandBox decoy set are smaller than the true active compounds of the target. In Tables 3 and 4, there is no clear trend due to the decoy set in the AUC and the *R* values. Thus, the idea of the UAP could be applied to the other decoy sets. However, many focused libraries have been developed: drug- and lead-like compound and fragment-like databases,

GPCR- (GPCR-focused) and kinase-oriented (kinase-focused) databases, databases for metallo protease, and others.^{37,38} These focused libraries include many drug-like compounds, thus we do not expect that the idea of the UAP would work very well in those cases.

An increase in the number of UAPs does not always increase the prediction accuracy (the correlation coefficient (R) between the AUC_{true} value and the AUC_{UAP}). The UAP_{avg} is a summation set of the UAP_{GPCR} , UAP_{DUD} , and UAP_{PDB} . The numbers of the compounds of the UAP_{GPCR} , UAP_{DUD} , and UAP_{PDB} are 73, 175, and 148, respectively. The number of compounds of the UAP_{avg} is 396. However, Table 3 shows that the R value obtained by the UAP_{PDB} is higher than that obtained by the UAP_{avg} . Our method may be available in the limited cases, and the prediction accuracy by the UAP should be also limited. If the active compounds of a protein that is similar to the target protein are available, then these active compounds should be used as the UAPs. Namely, most COX2 inhibitors are inhibitors of cyclooxygenase-1 (COX1).

5. CONCLUSION

We found that when a protein–compound docking program is used to predict the protein–compound affinity, the enzyme likely binds some drug-like compounds that are not its true ligands, in addition to its true ligand. Although this result may be an artifact, this finding is useful for improve in silico drug screening.

We introduced a drug-like compound set called the universal active probe (UAP). These compounds are the ligands of the GPCRs and the ligands selected from both the active compounds of the DUD and the protein–compound complex structures of the PDB. We performed structure-based in silico drug screening based on the ensemble of the target protein structures generated by the MD simulation. The hit ratio of true active compounds of a target protein correlated well with the hit ratio of the UAPs. Namely, the AUC of the true active compounds showed positive correlation to the AUC of the UAPs. Using this knowledge, we can assess the screening results obtained from the ensemble of target protein structures. There were some exceptions, but this assessment method is effective for many target proteins.

The idea of consensus scoring was also useful for improving the prediction accuracy by the UAPs. The average or minimum values of the AUCs of three different UAPs showed better correlations with the AUCs of the true active compounds than that of the single use of the AUC of UAPs.

ACKNOWLEDGMENT

This work was supported by grants from the New Energy and Industrial Technology Development Organization of Japan (NEDO) and the Ministry of Economy, Trade, and Industry (METI) of Japan.

APPENDIX A: MULTIPLE TARGET SCREENING (MTS) METHOD

We used a structure-based drug screening method that uses a protein–compound affinity matrix, called the MTS method.^{15,16} This is also a sort of “affinity fingerprint” approach. The basic idea of the MTS method is that the potential active compounds are those that show the strongest affinities to the target protein.

These compounds are sorted according to the docking scores. Thus, based on the protein–compound affinity matrix, these compounds are selected as the hit compounds. The protein set consists of 180 proteins listed in Appendix B, which were also used in our previous study.^{16,25} To perform the docking simulation, the SievGene/myPresto protein–compound docking program was used.¹¹ The docking program, the MTS screening tools, and the 3D structures of the used proteins are available on the web site (http://presto.protein.osaka-u.ac.jp/myPresto4/index_e.html).

APPENDIX B: PROTEIN SET USED.

The protein databank (PDB) identifier list of the basic protein set is: 1a28, 1a42, 1a4g, 1a4q, 1abe1, 1abe2, 1abf1, 1abf2, 1aco, 1ai5, 1aoe, 1apt, 1apu, 1aqw, 1atl, 1b58, 1b9v, 1bma, 1byb, 1byg, 1c1e, 1c5c, 1c83, 1cbs, 1cbx, 1cdg, 1ckp, 1com, 1coy, 1cps, 1cvu, 1d0l, 1d3h, 1dd7, 1dg5, 1dhf, 1dog, 1dr1, 1ebg, 1eed, 1eju, 1epb, 1epo, 1ets, 1f0r, 1f0s, 1f3d, 1fen, 1fkg, 1fki, 1fl3, 1glp, 1hdc, 1hfc, 1hos, 1hpb, 1hsb, 1hsl, 1htf1, 1htf2, 1hyt, 1ida, 1ivb, 1jap, 1lah, 1lcp, 1lic, 1lna, 1lst, 1mdr, 1mld, 1mmq, 1mrg, 1mts, 1mup, 1nco, 1ngp, 1nis, 1okl, 1pbd, 1phd, 1phg, 1poc, 1ppc, 1pph, 1pso, 1qbr, 1qbu, 1qpp, 1rds, 1rne, 1rnt, 1rob, 1snc, 1srj, 1tlp, 1tmn, 1tng, 1tnh, 1tni, 1tnl, 1tyl, 1xid, 1xie, 1yee, 2aad, 2ack, 2ada, 2cht, 2cmd, 2cpp, 2ctc, 2fox, 2gbp, 2ifb, 2pk4, 2qwk, 2tmn, 3cla, 3cpa, 3erd, 3ert, 3tpi, 4lbd, 4phv, 5abp1, 5abp2, 5cpp, 5er1, 6rnt, and 7tim. For 1abe, 1abf, 5abp, and 1htf, two protein pockets were prepared, since these proteins each bind two kinds of ligands.

Supporting Information Available: Protein structures and their inhibitors and schematic representation of the procedure used in the current study are summarized. This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- (1) Rueda, M.; Bottegoni, G.; Abagyan, R. Recipes for the selection of experimental protein conformations for virtual screening. *J. Chem. Inf. Model.* **2010**, *50*, 186–193.
- (2) Rueda, M.; Bottegoni, G.; Abagyan, R. Consistent improvement of cross-docking results using binding site ensembles generated with elastic network normal modes. *J. Chem. Inf. Model.* **2009**, *49*, 716–725.
- (3) Omagari, K.; Mitomo, D.; Kubota, S.; Nakamura, H.; Fukunishi, Y. A method to enhance the hit ratio by a combination of structure-based drug screening and ligand-based screening. *Adv. Appl. Bioinf. Chem.* **2008**, *1*, 19–28.
- (4) Warren, G. L.; Webster Andrews, C.; Capelli, A. M.; Clarke, B.; LaLonde, J.; Lambert, M. H.; Lindvall, M.; Nevins, N.; Semus, S. F.; Senger, S.; Tedesco, G.; Wall, I. D.; Woolven, J. M.; Peishoff, C. E.; Head, M. S. A critical assessment of docking programs and scoring functions. *J. Med. Chem.* **2006**, *49*, 5912–5931.
- (5) Kontoyianni, M.; Sokol, G. S.; McClellan, L. M. Evaluation of library ranking efficacy in virtual screening. *J. Comput. Chem.* **2005**, *26*, 11–22.
- (6) McGovern, S. L.; Shoichet, B. K. Information decay in molecular docking screens against holo, apo, and modeled conformations of enzymes. *J. Med. Chem.* **2003**, *46*, 2895–2907.
- (7) Diller, D. J.; Li, R. Kinase, homology models, and high throughput docking. *J. Med. Chem.* **2003**, *46*, 4638–4647.
- (8) Bissantz, C.; Bernard, P.; Hibert, M.; Rognan, D. Protein-based virtual screening of chemical database. II. Are homology models of G-protein coupled receptors suitable targets. *Proteins* **2003**, *50*, 5–25.
- (9) Oshiro, C.; Bradley, E. K.; Eksterowicz, J.; Evensen, E.; Lamb, M. L.; Lancot, K.; Putta, S.; Stanton, R.; Grootenhuys, D. J. Performance of 3D-database molecular docking studies into homology models. *J. Med. Chem.* **2004**, *47*, 764–767.
- (10) DeWesse-Scott, C.; Moul, J. Molecular modeling of protein function regions. *Proteins* **2004**, *55*, 942–961.

- (11) Rueda, M.; Bottegioni, G.; Abagyan, R. Recipes for the selection of experimental protein conformations for virtual screening. *J. Chem. Inf. Model.* **2010**, *50*, 186–193.
- (12) Barakat, K.; Mane, J.; Friesen, D.; Tuszyński, J. Ensemble-based virtual screening reveals dual-inhibitors for the p53-MDM2/MDMX interactions. *J. Mol. Graphics Modell.* **2010**, *28*, 555–568.
- (13) Craig, I. R.; Essex, J.; Spiegel, K. Ensemble docking into multiple crystallographically derived protein structures: an evaluation based on the statistical analysis of enrichment. *J. Chem. Inf. Model.* **2010**, in press.
- (14) Fukunishi, Y.; Mikami, Y.; Nakamura, H. Similarities among receptor pockets and among compounds: Analysis and application to in silico ligand screening. *J. Mol. Graphics Modell.* **2005**, *24*, 34–45.
- (15) Fukunishi, Y.; Mikami, Y.; Kubota, S.; Nakamura, H. Multiple target screening method for robust and accurate in silico ligand screening. *J. Mol. Graphics Modell.* **2005**, *25*, 61–70.
- (16) Fukunishi, Y.; Kubota, S.; Nakamura, H. Noise reduction method for molecular interaction energy: application to in silico drug screening and in silico target protein screening. *J. Chem. Inf. Model.* **2006**, *46*, 2071–2084.
- (17) Huang, N.; Shoichet, B. K.; Irwin, J. J. Benchmarking sets for molecular docking. *J. Med. Chem.* **2006**, *49*, 6789–6801.
- (18) Fukunishi, Y.; Kubota, S.; Nakamura, H. Finding ligands for G-protein coupled receptors based on the protein-compound affinity matrix. *J. Mol. Graphics Modell.* **2007**, *25*, 633–43.
- (19) Pipeline Pilot; Accelrys, Inc.: San Diego, CA; <http://accelrys.com/products/pipeline-pilot/>.
- (20) Hert, J.; Willett, P.; Wilton, D. J.; Acklin, P.; Azzaoui, K.; Jacoby, E.; Schuffenhauer, A. Comparison of topological descriptors for similarity-based virtual screening using multiple bioactive reference structures. *Org. Biomol. Chem.* **2004**, *2*, 3256–3266.
- (21) Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. Development and testing of a general amber force field. *J. Comput. Chem.* **2004**, *25*, 1157–1174.
- (22) Fukunishi, Y.; Mikami, Y.; Nakamura, H. The filling potential method: A method for estimating the free energy surface for protein-ligand docking. *J. Phys. Chem. B* **2003**, *107*, 13201–13210.
- (23) Gasteiger, J.; Marsili, M. Iterative partial equalization of orbital electronegativity — a rapid access to atomic charges. *Tetrahedron* **1980**, *36*, 3219–3228.
- (24) Gasteiger, J.; Marsili, M. A new model for calculating atomic charges in molecules. *Tetrahedron Lett.* **1978**, 3181–3184.
- (25) Fukunishi, Y.; Mashimo, T.; Orita, M.; Ohno, K.; Nakamura, H. In silico fragment screening by replica generation (FSRG) method for fragment-based drug design. *J. Chem. Inf. Model.* **2009**, *49*, 925–933.
- (26) Fukunishi, Y.; Nakamura, H. Improvement of protein–compound docking scores by using amino-acid sequence similarities of proteins. *J. Chem. Inf. Model.* **2008**, *48*, 148–156.
- (27) Nissink, J. W. M.; Murray, C.; Hartshorn, M.; Verdonk, M. L.; Cole, J. C.; Taylor, R. A new test set for validating predictions of protein–ligand interaction. *Proteins* **2002**, *49*, 457–471.
- (28) Case, D. A.; Darden, T. A.; Cheatham, T. E., III.; Simmerling, C. L.; Wang, J.; Duke, R. E.; Luo, R.; Merz, K. M.; Wang, B.; Pearlman, D. A.; Crowley, M.; Brozell, S.; Tsui, V.; Gohlke, H.; Mongan, J.; Hornak, V.; Cui, G.; Beroza, P.; Schafmeister, C.; Caldwell, J. W.; Ross, W. S.; Kollman, P. A. *AMBER 8*; University of California: San Francisco, CA, 2004.
- (29) Jorgensen, W. L.; Chandrasekhar, J.; Madura, J. D.; Impey, R. W.; Klein, M. L. Comparison of simple potential functions for simulating lipid water. *J. Chem. Phys.* **1983**, *79*, 926–935.
- (30) Ryckaert, J. P.; Ciccotti, G.; Berendsen, H. J. C. Numerical integration of the cartesian equations of motion of a system with constraints: Molecular dynamics of n-alkanes. *J. Comput. Phys.* **1977**, *23*, 327–341.
- (31) Greengard, L.; Rokhlin, V. A fast algorithm for particle simulations. *J. Comput. Phys.* **1987**, *73*, 325–348.
- (32) Fukunishi, Y.; Sugihara, Y.; Mikami, Y.; Sakai, K.; Kusudo, H.; Nakamura, H. Advanced in-silico drug screening to achieve high hit ratio-development of 3D-compound database. *Synthesiology* **2009**, *2*, 60–68.
- (33) Fukunishi, Y.; Nakamura, H. Prediction of protein-ligand complex by docking software guided by other complex structures. *J. Mol. Graphics Modell.* **2008**, *26*, 1030–1033.
- (34) Fukunishi, Y.; Nakamura, H. A new method for in-silico drug screening and similarity search using molecular-dynamics maximum-volume overlap (MD-MVO) method. *J. Mol. Graphics Modell.* **2009**, *27*, 628–636.
- (35) Coleman, R. G.; Sharp, K. A. Protein pockets: inventory, shape, and comparison. *J. Chem. Inf. Model.* **2010**, in press.
- (36) Fujii, I.; Sugaya, N.; Nakano, T.; Hasegawa, S.; Yamamoto, M.; Kaminuma, T.; Hirayama, N. Essential chemical characteristics for drugs. *Chem-Bio Inf. J.* **2001**, *1*, 18–22.
- (37) Willett, P.; Barnard, J. M.; Downs, G. M. Chemical similarity searching. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 983–996.
- (38) Jimonet, P.; Jager, R. Strategies for designing GPCR-focused libraries and screening sets. *Curr. Opin. Drug Discovery Devel.* **2004**, *7*, 325–333.

CI100108P