# Similarities among receptor pockets and among compounds: Analysis and application to in silico ligand screening

Yoshifumi Fukunishi [a,*], Yoshiaki Mikami [b], Haruki Nakamura [a,c]

[a] Biological Information Research Center (BIRC), National Institute of Advanced Industrial Science and Technology (AIST), 2-41-6, Aomi, Koto-ku, Tokyo 135-0064, Japan
[b] Japan Biological Information Research Center (JBIRC), Japan Biological Informatics Consortium (JBIC), 2-41-6, Aomi, Koto-ku, Tokyo 135-0064, Japan
[c] Institute for Protein Research, Osaka University, 3-2 Yamadaoka, Suita, Osaka 565-0871, Japan

## Abstract

We developed a new method to evaluate the distances and similarities between receptor pockets or chemical compounds based on a multi-receptor versus multi-ligand docking affinity matrix. The receptors were classified by a cluster analysis based on calculations of the distance between receptor pockets. A set of low homologous receptors that bind a similar compound could be classified into one cluster. Based on this line of reasoning, we proposed a new in silico screening method. According to this method, compounds in a database were docked to multiple targets. The new docking score was a slightly modified version of the multiple active site correction (MASC) score. Receptors that were at a set distance from the target receptor were not included in the analysis, and the modified MASC scores were calculated for the selected receptors. The choice of the receptors is important to achieve a good screening result, and our clustering of receptors is useful to this purpose. This method was applied to the analysis of a set of 132 receptors and 132 compounds, and the results demonstrated that this method achieves a high hit ratio, as compared to that of a uniform sampling, using a receptor–ligand docking program, Sievgene, which was newly developed with a good docking performance yielding 50.8% of the reconstructed complexes at a distance of less than 2 Å RMSD.
© 2005 Elsevier Inc. All rights reserved.

*Keywords:* Receptor–ligand docking; Flexible docking; Docking score; Cluster analysis; Receptor classification; Database enrichment

## 1. Introduction

For structure-based drug design, similarities between receptor pockets are an important aspect to investigate in terms of the selectivity of chemical compounds and for the identification of new candidate drugs. An active and useful compound will be able to bind to the target pocket as well as to other, similar pockets. In the latter case, the similar pocket can be associated with side effects. Recent molecular-targeting drugs require high selectivity in order to avoid such side effects, and therefore it is necessary to identify the receptor pockets that are most similar to the target pocket. Then, the similarities and differences between a target

pocket and similar pockets must be considered. For example, in order to develop non-steroidal anti-inflammatory drugs (NSAIDs), it is necessary to consider the 3D structure of cyclooxygenase (COX)-2, which is the target, and that of COX-1, which similar to the target pocket of COX-2 [1–3]. The identification of similarities between receptor pockets is useful in the search for new drugs. If a known drug is similar to a compound that binds to a target pocket, that drug is likely to bind to the target pocket as well. Inhibitors of HIV protease have been developed based on inhibitors of renin, a protein that is homologous to HIV protease [4].

Several methods have been developed to evaluate the similarity between receptor pockets [5–14]. The most widely used method involves determination of the sequence homology between amino-acid sequences of proteins. This method has the advantages of being both

* Corresponding author. Tel.: +81 3 3599 8290; fax: +81 3 3599 8099.
*E-mail address:* y-fukunishi@jbirc.aist.go.jp (Y. Fukunishi).

quite rapid and powerful, although a number of exceptions have been reported. For example, there are a number of ATP binding proteins exhibiting quite low homology with each other [15]. Recently, the direct comparison of 3D protein surfaces has been made available [5–7]. Such methods can reveal the degree of similarity between two protein surfaces, and such approaches can also point out the differences between receptor pockets in quantitative terms. In particular, these recently developed methods have been used to successfully identify pairs of non-homologous proteins that bind to the same compound. Once it is possible to determine the difference and/or similarity between two pockets, we can apply both a principal component analysis and a cluster analysis to a given set of pockets in order to classify them.

The search for similarities among chemical compounds is common in the field of pharmaceutical science. Similar compounds are expected to bind to the same target pocket; thus, the identification of similar compounds is a basic technique used in drug screening, and a number of methods have been developed for this purpose. For example, the CATS descriptor represents the spatial distribution of functional groups in a compound [16]. The BCUT descriptor is a set of several numbers that represents the chemical structure of the compound, and the BCUT reveals the diversity of a chemical compound library [17,18]. Similarities between functional groups are determined based on a propensity map, which represents the spatial distribution of other functional groups around the functional group in question [19,20]. The difference or similarity between two functional groups is determined as the difference between the two respective propensity maps. However, it has remained necessary to develop a method for the direct comparison of receptor–ligand interactions.

Thus, we developed a novel method for the evaluation of the distance and similarity between protein (receptor) pockets and chemical compounds; this method is based on the receptor–ligand docking interaction. To this end, we prepared two datasets, one being a set of receptor pockets, and the other a set of ligand chemical compounds. Then, the protein–ligand interaction panel $s(i, k)$ is calculated for the $i$-th pocket and the $k$-th ligand, where $i$ and $k$ assume the serial numbers for all pockets and ligands, respectively. The distance and similarity are defined using the following panel: $s(i, k)$. We also proposed a novel in silico screening scheme using a number of targets based on the panel $s(i, k)$. The combination of this scheme, used together with a modification of the recently developed Multiple Active Site Correction (MASC) score, was found to give a good hit ratio, as compared to that of a random sampling [21].

In order to carry out the screening procedure, a new protein–ligand docking program, designated as "sievgene", was developed, although many protein–ligand docking programs have been reported [22–31]. One of the serious problems in protein–ligand docking is the occurrence of poor atomic contact due to rough docking, which originates from a van der Waals repulsion between the protein and the ligand. The scoring function of our method is based on a rough shape of the protein surface to reduce structural noises. The conventional potential function is applied to the outer region of the protein, and a smooth virtual function is applied to the inner region of the protein. This score function reduced the artifacts due to atomic contact, and it correctly predicted $\sim$50% of the complex conformations.

## 2. Methods

### 2.1. Definition of distances and similarities among receptor pockets

The distance between two protein pockets is determined based on the protein–ligand interaction matrix. Moreover, it is also possible to determine the similarity between two pockets, as well as the relative selectivity of pockets. The same discussion as that given here can be applied to the compounds.

Here, we will prepare a set of pockets $P = \{p_1, p_2, p_3, \ldots p_M\}$, where $p_i$ represents the $i$-th pocket. The total number of pockets is $M$. We will also prepare a set of compounds $X = \{x^1, x^2, \ldots x^N\}$, where $x^k$ represents the $k$-th compound. The total number of compounds is $N$. For each pocket $p_i$, all compounds of the set $X$ are docked to the pocket $p_i$, and $v_i$ represents the score vector of $p_i$. Then, $s_i^k$ represents the binding score between the $i$-th pocket and the $k$-th compound. Here, $s_i^k$ corresponds to the absolute value of the free energy of binding, and $s_i^k \geq 0$

$$v_1 = (s_1^1, s_1^2, s_1^3, \ldots, s_1^N)$$

$$v_i = (s_i^1, s_i^2, s_i^3, \ldots, s_i^N)$$

$$v_M = (s_M^1, s_M^2, s_M^3, \ldots, s_M^N)$$

The distance between the $i$-th pocket and the $j$-th pocket is defined as follows

$$D_{ij}^P = \sqrt{(v_i - v_j)^2} = \sqrt{\sum_{k=1}^{N} (s_i^k - s_j^k)^2} \tag{1}$$

or

$$D_{ij}^P = \sum_{k=1}^{N} \left| s_i^k - s_j^k \right| \tag{2}$$

$D_{ij}^P$ satisfies the following three conditions, which are sufficient for a definition of distance, $D_{ii}^P = 0$, $D_{ij}^P = D_{ji}^P$, and $D_{ih}^P + D_{hj}^P \geq D_{ij}^P$. Eqs. (1) and (2) indicate that if two pockets bind the same compounds with the same binding free energy, the distance between the two pockets is zero. If two pockets

are unable to bind any compound, the distance between them becomes zero.

Also, the similarity between the $i$-th pocket and the $j$-th pocket is defined as follows

$$S_{ij}^P = \frac{v_i \times v_j}{|v_i||v_i|} \tag{3}$$

Here the dot means inner product and $S_{ij}^P \geq 0$, because all $s_i^k \geq 0$, and $S_{ij}^P \leq 1$, based on the given definition. If two pockets bind the same compounds with the same binding free energy, the similarity between them is 1. If the pockets cannot bind any of the compounds tested, we cannot determine the degree of similarity between them. Considering both the distance and the similarity together, we can apply a cluster analysis and a principal component analysis to the set of pockets.

As mentioned above, we can apply the same discussion to the compounds. The score vector $w^k$ of the $k$-th compound can be determined as follows

$$w^k = (s_1^k, s_2^k, s_3^k, \ldots, s_M^k)$$

The distance between the $k$-th compound and the $l$-th compound is defined as follows

$$D_{kl}^C = \sqrt{(w^k - w^l)^2} = \sqrt{\sum_{i=1}^{M}(s_i^k - s_i^l)^2} \tag{4}$$

or

$$D_{kl}^C = \sum_{i=1}^{M}\left|s_i^k - s_i^l\right| \tag{5}$$

$D_{kl}^C$ sufficiently satisfies the conditions for distance, as does $D_{ij}^P$. The similarity between compounds $k$ and $l$ is thus defined as

$$S_{kl}^C = \frac{w^k \times w^l}{|w^k||w^l|} \tag{6}$$

Here, $S_{kl}^C$ satisfies the same condition as $S_{ij}^P$.

### 2.2. Docking method

The concepts of our docking method are a meta-heuristic minimum search and a robust scoring function. The meta-heuristic approach is a combination of a global minimum search and a heuristic local minimum search, for example the steepest descent method. The global search roughly screens the overall potential surface to find candidates for the potential global minimum. The successive local minimum search finds the exact minimum in the local potential basin. As regards the global search method, we superimposed ligand atoms to the binding site of the receptor using the geometric hashing method [32]. As regards the local minimum search method, we adopted the steepest descent method. The details of this approach are as follows

- *Step* 1 The pocket is indicated by a set of reference points, which are the atom coordinates of the ligand of a crystalline protein–ligand complex.
- *Step* 2 Electrostatic potential field on the accessible surface of the receptor is calculated to find total 30 potential minima and maxima, also 30 points, on which the potential is nearly zero, are found. Triangles are generated to connect these points; the data regarding these triangles are recorded in a hash table.
- *Step* 3 The program reads the coordinates and the atom types of the ligand molecule and then generates its conformers. The dihedral angles are randomly incremented every $120°$.
- *Step* 4 The global search program chooses any three atoms of the ligand molecule and superimposes the ligand molecule onto the receptor surface according to the geometric hash method, and the energy score is evaluated. The results are sorted according to the order of the energy scores.
- *Step* 5 The next global search is reinitiated with different conformers. The new search area is limited to the region near the predicted ligand coordinates described in Step 4. The hash table is reconstructed for the new search area.
- *Step* 6 The ligand-conformer generation program generates the new conformers, which are similar to certain ligand structures of candidate ligand–protein complexes. The similarity is then measured by the RMSD of all atoms; in this study, RMSD threshold was set at $<3.0$ Å.
- *Step* 7 Returning to Step 4, a global search is conducted based on the new search area and the new conformers. This process is repeated five times.
- *Step* 8 Starting from the initial coordinates the system reaches the optimal complex structure using the steepest descent method with the Grid potential of the receptor force field. The AMBER-type molecular force field is used. The equilibrium bond lengths and bond angles are set to the initial coordinates of the ligand molecule. The bond and torsion force constants are assigned to arbitrary default values.

### 2.3. Scoring function

The receptor–ligand interactions accounted for by our method are van der Waals (vdW), Coulomb, hydrogen bond, and hydrophobic interactions. A grid potential represents these interactions, assuming that the receptor structure is rigid. The rigid receptor model cannot represent induced fitting, and therefore it is possible to overestimate the repulsion energy due to atomic contact. To reduce this structural noise, the repulsive part of the ligand–receptor potential is modified in order to underestimate the atomic contact.

Using our method, the space is divided into two regions, i.e., the inner and outer regions of the protein, according to the accessible surface of the receptor. In the inner region of the receptor, we apply an artificial smoothing function, and the force field of the outer and the inner regions are smoothly

connected at the interface, which is the accessible surface of the receptor.

The potential function for van der Waals interactions in the outer region is

$$E_{vdW} = w(r) \sum_{a,b} 4\varepsilon_{a,b}\left(\left(\frac{\sigma_{ab}}{R_{ab}}\right)^{12} - \left(\frac{\sigma_{ab}}{R_{ab}}\right)^6\right) \quad (7)$$

where $R_{ab}$ is the distance between the a-th atom of the ligand and the b-th atom of the receptor, and $w$ is a weight function

$$w(r) = 1 - e^{-C_1 r} \quad (8)$$

where $C_1$ is a constant and $r$ is the minimum distance from the a-th atom to the interface, which separates the inner and the outer regions. In this study, $C_1$ was set to 2.0 Å$^{-1}$. In the inner regions, the potential function for vdW interactions is

$$E_{vdW} = Dr \quad (9)$$

where $D$ is a constant and $r$ is the minimum distance from the interface. In this study, $D$ was set to 5.0 kcal/mol/Å. At the interface, the values of both Eqs. (8) and (9) become zero, such that the vdW potential surface is continuous.

The potential function for Coulomb interactions is given by Eq. (10) for the protein's outer region, and it is given by Eq. (11) for the inner region of the protein. Thus

$$E_{elec} = \sum_{a,b} \frac{332 q_a q_b}{\varepsilon \times r_{ab}^2} \quad (10)$$

where $q_a$, $q_b$, and $r_{ab}$ are the atomic charges of the a-th atom of the ligand and the b-th atom of the receptor, and the distance between the a-th and the b-th atoms, respectively. $\varepsilon$ is a dielectric constant, and we use a uniform value, $\varepsilon = 4.0$, in this case

$$E_{elec} = w(r_{ASA}) \sum_{a,b} \frac{332 q_a q_b}{\varepsilon \times r_{ASA}^2} \quad (11)$$

where $r_{ASA}$ is the minimum distance between the a-th atom of the ligand to the receptor-accessible surface. In addition, $w$ is the following weight function

$$w(R) = \frac{1}{1 + R} \quad (12)$$

At the interface, the value of Eq. (11) becomes the same as that given by Eq. (12).

The potential function for a hydrogen bond is

$$E_{H-bond} = \sum_{a,b} C_2 \, e^{-((R_{ab} - d_{OH})/r_1)^2 - ((\theta - \theta_0)/r_2)^2} \quad (13)$$

where $R_{ab}$ and $\theta$ are the distance and angle between the a-th hydrogen donor of the ligand and the b-th hydrogen acceptor of the receptor, or vice versa. Here, $C_2$, $r_1$, $d_{OH}$, $r_2$ and $\theta_0$ are constant and they are set to $-3.0$ kcal/mol, 3.0 Å, 1.2 Å, 40.0°, and 180.0°, respectively. This potential function is similar to that of PRO_LEADS, which successfully reproduced the binding affinity of various protein–ligand complexes [26].

The potential function for hydrophobic interaction is

$$E_{ASA} = \begin{cases} \sum_{a,b} f(A + C_3 R_{ab}) : R_{ab} \leq \sigma_a + \sigma_b + 2r_{prob} \\ 0 : R_{ab} > \sigma_a + \sigma_b + 2r_{prob} \end{cases}$$

$$(14)$$

where $A = -4\pi\{(\sigma_a + r_{prob})^2 + (\sigma_b + r_{prob})^2\}$ and $C_3 = -A/(\sigma_a + \sigma_b + 2r_{prob})$. Here, $\sigma_a$ and $r_{prob}$ are the vdW radius of the a-th atom and the probe radius. In this study, $f$ is an atomic solvation parameter and it is set to 10.0 cal/mol/Å$^2$ for all atoms. Although the pairwise potential in Eq. (14) is only a rough approximation for the hydrophobic interaction, pairwise potentials have been used frequently in docking programs and have provided a good estimate [28,33].

We can numerically smooth the grid potential in order to reduce the structural noise

$$v_{\alpha,\beta,\gamma} = \frac{n \times v_{\alpha,\beta,\gamma} + \sum_{\alpha'=\alpha-1,\alpha+1,\beta'=\beta-1,\beta+1,\gamma'=\gamma-1,\gamma+1} v_{\alpha',\beta',\gamma'}}{n + 6} \quad (15)$$

where $v_{\alpha,\beta,\gamma}$ is a score value at grid point $(\alpha, \beta, \gamma)$, and $n$ is set to 2. This smoothing process is then applied six times.

Finally, for the actual potential score calculations, we can apply the Lagrange first-order interpolation to the grid potential. The dimensionless raw docking score that is used for receptor–ligand docking is

$$S_{raw} = g(E_{vdW} + E_{elec} + E_{H-bond} + E_{ASA}) \quad (16)$$

where $g$ is a parameter and it is set to 0.01 mol/kcal.

We can also calculate the MASC score based on the raw docking score. The MASC score $S'_{ik}$ for the i-th pocket and the k-th compound has been reported by Vigers and Rizzi as follows [21]

$$S'_{ik} = \frac{(S_{ik} - \mu_k)}{\sigma_k} \quad (17)$$

where $S_{ik}$ is the raw docking score, and $\mu_k$ and $\sigma_k$ are the average and the standard deviation of the raw docking score across all pockets for the k-th compound, respectively.

Following the above docking scheme, a new program, sievgene, was developed to provide the scores.

### 2.4. Screening method

We proposed a new in silico screening method by which the compounds in the database are docked to multiple targets that could function as representative receptor pockets. The new score obtained by this method is a slightly modified version of the multiple active site correction (MASC) score. To apply Eq. (17), the distribution of docking scores must approach a Gaussian distribution. In general, extreme values must be removed from the data in order to improve the distribution. For example, a positive score corresponding to a positive binding free energy value, is meaningless by the definition of binding free energy, and such data must be

excluded from the analysis. Thus, only raw docking scores of less than $S^0$ can be used for the analysis. In this study, $S^0$ was set to zero. Moreover, if the raw docking score was exactly equal to the binding free energy, the raw docking score was more accurate than the MASC score. Therefore, the accuracy of the MASC score depends on the quality of the raw docking score. The new modified MASC score ($S_{new}$) is a linear combination of the MASC score ($S_{MASC}$) and the raw docking score ($S_{raw}$). Here, $\lambda$ is a parameter

$$S_{new} = \lambda \times S_{raw} + (1 - \lambda)S_{MASC} \qquad (18)$$

Receptors that are distant from the target receptor are excluded from the analysis, and the modified MASC scores, $S_{new}$, are calculated for the selected receptors. Let us assume that the receptor dataset consists of two receptor groups. The receptors of one group show high docking scores for almost all compounds, and the receptors of the other group show low docking scores for almost all compounds. In this case, the receptors of the former group form one cluster, and the receptors of the latter group form another cluster. In addition, the distribution of the docking scores is different from a Gaussian distribution. Thus, the distances of each receptor between the target receptors are sorted, and the top $R_a$ % of the total receptors is used for the analysis, according to the conventional screening method, used in concert with the modified MASC score. We refer to this method as the "receptor selection (RS) method".

## 3. Preparation of materials

To evaluate our docking method, we performed a protein–ligand docking simulation based on the known complex structures registered in the Protein Data Bank. Here, 132 complexes were selected from the database, which was used in the evaluation of the GOLD and FlexX [34]. This data set contains a rich variety of proteins and compounds whose structures were all determined by high quality experiments with a resolution less than 2.5 Å. The lack of atom coordinates is almost zero and the all-atomic structures around the ligand pockets are quite reliable. Thus, this data set was used in the clustering analysis of proteins, compounds and in silico screening. We removed from the original data set those complexes containing a covalent bond between the protein and ligand. The PDB identifiers are summarized in Appendix A. All water molecules and cofactors were removed from the proteins and all missing hydrogen atoms were added to form the all atom models of proteins. The conformations of all ligands were randomized before the docking simulation.

The size distribution of ligands was as follows: ratio of 1–9 atoms, 3.6%; ratio of 10–19 atoms, 15.2%; ratio of 20–29 atoms, 30.9%; ratio of 30–39 atoms, 15.4%; ratio of 40–49 atoms, 15.8%; ratio of 50–59 atoms, 10.6%; and the ratio of more than 60 atoms was 16.1%. The average ligand size was 37.1 atoms.

The atomic charges of each ligand were determined by the restricted electrostatic point charge (RESP) procedure using HF/6-31G*-level quantum chemical calculations [35]. We used GAMESS and Gaussian98 to perform the quantum chemical calculations [36,37]. The atomic charges of the proteins were the same as the atomic charges in AMBER parm99 [35,38].

## 4. Results and discussion

### 4.1. Docking accuracy and classification results

Our docking program, sievgene, reconstructed 18.9, 50.8, and 59.8% of a total of 132 complexes with RMSD values of <1, 2, and 3 Å, respectively, where the RMSD was calculated between all atom positions, with the exception of the H atoms of each docked compound and the corresponding atom positions in the complex crystal structure. Using almost the same dataset, DOCK [22], FlexX [23], and GOLD [24] were reported to reconstruct the complexes of 39, 51, and 56%, with RMSDs <2 Å, respectively [25]. Thus, the accuracy of our docking program was found to be as good as that of the popular docking programs. The average CPU time was 845 s on a single CPU of a COMPAQ ALPHA ES45. For eleven particular compounds, the computation took more than one hour. With the exception of these compounds, the average CPU time was 542 s.

Sievgene exhibited high prediction accuracy and calculation speed when small ligands were analyzed. It reconstructed 20.3, 59.4, and 65.2% of a total of 69 complexes in which the number of ligand atoms was less than 40, with RMSD values of <1, 2, and 3 Å, respectively, and the average computational time was 179 CPU seconds. Sievgene reconstructed 17.5, 41.3, and 52.4% of a total of 63 complexes in which the number of ligand atoms was greater than 40, with RMSD <1, 2, and 3 Å, respectively, and the average computational time was 960 CPU seconds. Sievgene was found to randomly generate several hundreds of ligand conformers, without taking the conformational energy into account; whereas the conventional docking programs also evaluate the conformational energy of ligand. The conformation generation scheme remains to be improved in the future.

Cluster analysis was applied to all 132 receptors, based on the receptor–compound score matrix. Here, 132 proteins were divided at cluster level 4.45, and these proteins were divided into seven clusters (cluster 1–cluster 7); it should be noted that it was possible to divide the proteins into any number of clusters. There were several singleton clusters, and these small clusters were merged into the nearest cluster.

The clustering results are summarized in Table 1. Cluster 1 consisted primarily of sugar-binding proteins, and all seven sugar-binding proteins in the database were included in cluster 1. Cluster 2 consisted primarily of serine proteases. Cluster 3 consisted primarily of neuraminidases and hydrolases. All five neuraminidases in the database were

Table 1
Cluster classification results

| Cluster | Proteins | | | | | Feature | Content in cluster[a] | Content in database[b] |
|---|---|---|---|---|---|---|---|---|
| Cluster 1 | 1abe1 | 1abe2 | 5app1 | 5app2 | 1abf1 | Sugar-binding | 70%(7) | 5.3%(7) |
| | 1abf2 | 2gbp | 1lst | 1lah | 1ebg | | | |
| Cluster 2 | 1tni | 1tng | 1tnl | 1tnh | 1f0s | Serine protease | 20%(4) | 8.3%(11) |
| | 1hfc | 1atl | 1f0r | 1nqp | 1mrg | | | |
| | 1xid | 1hyt | 1f3d | 1xie | 1ai5 | | | |
| | 2ack | 3erd | 1a28 | 2ada | 1dog | | | |
| Cluster 3 | 1b9v | 1a4q | 2qwk | 1a4g | 1a42 | Neuraminidase | 26.3%(5) | 3.8%(5) |
| | 1ejn | 2tmn | 1snc | 1ivb | 1hsb | Hydrolase | 26.3%(5) | 12.8%(17) |
| | 1aqw | 1glp | 3tpl | 1fl3 | 1mdr | | | |
| | 1cps | 1cbx | 1pbd | 1hsl | | | | |
| Cluster 4 | 1tlp | 1lna | 1tmn | 5er1 | 1rne | HIV protease | 19%(4) | 4.5%(6) |
| | 1pso | 1ets | 1gbr | 1htf1 | 1htf2 | Acid protease | 28.6%(6) | 4.5%(6) |
| | 1byg | 1byb | 1hos | 1hdc | 1dd7 | | | |
| | 1ida | 1epo | 1apt | 1eed | 1apu | | | |
| | 2ctc | | | | | | | |
| Cluster 5 | 1pph | 1mts | 1ppc | 3cla | 1d0l | Endonuclease | 16.0%(4) | 3.0%(4) |
| | 1srj | 1rob | 1mmq | 1jap | 2aad | Serine protease | 24.0%(6) | 8.3%(11) |
| | 1rnt | 1fki | 4est | 1bma | 2pk4 | | | |
| | 2fox | 1mup | 6rnt | 1tyl | 1nco | | | |
| | 1rds | 1cdg | 1fkg | 1nis | 1aco | | | |
| Cluster 6 | 5cpp | 1phd | 2cpp | 1png | 1dr1 | Oxidoreductase | 45.8%(11) | 12.1%(16) |
| | 1coy | 1cvu | 3ert | 4lbd | 1dg5 | | | |
| | 1aoe | 1ckp | 1poc | 1lic | 1dhf | | | |
| | 1epb | 1cbs | 2ifb | 1fen | 1qbu | | | |
| | 1hpv | 4phv | 2cnt | 1d3h | | | | |
| Cluster 7 | 1com | 1c1e | 1okl | 1c5c | 1yee | Catalytic antibody | 23.1%(3) | 2.2%(3) |
| | 1b58 | 7tim | 1c83 | 3cpa | 1lcp | Oxidoreductase | 23.1%(3) | 12.1%(16) |
| | 1qpq | 2cmd | 1mld | | | | | |

[a] Numbers in parenthesis represent the number of proteins in each cluster.
[b] Numbers in parentheses represent the number of proteins in the entire database.

included in cluster 3. Cluster 4 consisted primarily of HIV proteases and acid proteases. All six acid proteases in the database were included in cluster 4. Cluster 5 consisted primarily of endonucleases and serine proteases. All four endonucleases in the database were included in cluster 5. Cluster 6 consisted primarily of oxidoreductases. Cluster 7 consisted primarily of catalytic antibodies and oxidoreductases. All three catalytic antibodies in the database were included in cluster 7. This result shows that our clustering method works well, even if only ~50% of complex structures were correctly predicted.

Four out of 11 serine proteases (1tni, 1tng, 1tnh, and 1tnl) are included in Cluster 2 and 6 out of 11 serine proteases (1bma, 4est, 1mts, 1ppc, 1pph, and 2pk4) are included in Cluster 5. Some of these proteins are equivalent to each other, and there is only one unique serine protease in Cluster 2. There are three unique serine proteases in Cluster 5, and the amino acid sequence identities among them are 33.3, 38.0, and 66.7%. This means that the most of serine proteases are included in Cluster 5. The sequence identities among the serine protease in Cluster 2 and the three-serine proteases in Cluster 5 are 33.3, 38.0 and 97.4%. The serine protease in Cluster 2 is quite similar to one of the serine

proteases (1mts, 1ppc, 1pph) in Cluster 5. Even if the sequences are the same, the 3D structures could be different from each other because of the induced fitting with different ligand. Thus, the structure change by induced fitting and the low docking accuracy may have been the cause of the misassignment in the cluster analysis.

Fig. 1a shows the partial receptor-classification results of 31 pockets versus 132 compounds out of 132 pockets versus 132 compounds. The definition of distance is given by Eq. (1). The receptors were arbitrarily classified into four clusters for convenience. Similar receptor pockets were classified in the same cluster. The sugar-binding proteins are in cluster A, i.e., 1abe1, 1abe2, 1abf1, 1abf2, 5abp1, and 5abp2 are the same proteins (L-arabinose-binding protein). 1abe1 is a complex with alpha-L-arabinose and 1abe2 is a complex with beta-L-arabinose. 1abf1 is a complex with alpha-D-fucose and 1abf2 is a complex with beta-D-fucose. 5abp1 is a complex with alpha-D-galactose and 5abp2 is a complex with beta-D-galactose. These proteins are in one cluster. Note that the amino-acid sequences of 1abe1, 1abe2, abf1, abf2, 5abp1, and 5abp2 are exactly the same. In addition, the amino-acid sequences of the acid proteases (1tlp, 1lna, and 1tmn; thermolysin) are identical.
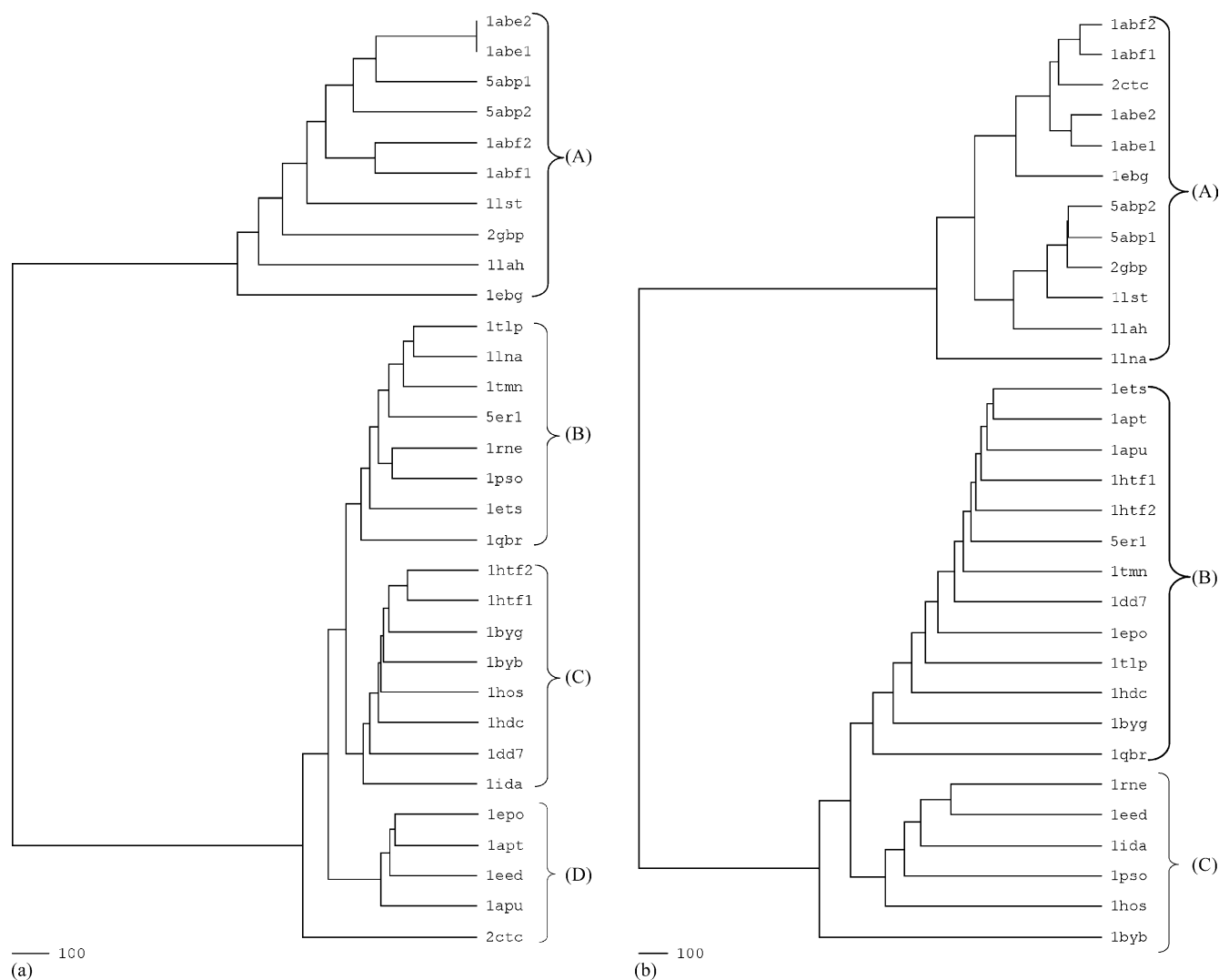
Fig. 1. (a) Cluster analysis of 31 receptor pockets from a total of 132. (b) Cluster analysis of 31 compounds from a total of 132. This phylogenetic tree was drawn with Tree View X [40].

Non-homologous proteins with similar ligands were classified into the same cluster. For example, 2gbp, which binds D-glucose, showed no homology (sequence identity = 24.7%) with 1abf, 1abe or 5abp, and these proteins were classified into the same cluster. Metalloproteases (1tlp, 1lna, and 1tmn) were included in cluster B, and HIV proteases (1htf1, 1htf2, 1hos, and 1ida) and oxidoreductases (1hdc and 1dd7) were included in cluster C. Four acid proteases (1epo, 1apt, 1eed, and 1apu) were included in cluster D, and two acid proteases (5er1 and 1rne) were included in cluster B. Homologous proteins were classified into the same cluster, i.e., in cluster C. HIV protease 1hos was chosen as a reference sequence, and the sequence identities of 1htf1, 1htf2, and 1ida were 100, 100, and 48.48%, respectively. In cluster D, the sequences of 1epo and 1eed were equivalent. In addition, the sequences of 1apt and 1apu were equivalent, but the sequence identity between 1epo and 1apt was 53.9%. Moreover, the neuraminidases (1ivb, 1a4g, 1a4q, 1b9v, and 2qwk) were classified in cluster

3, whereas the sequence identity between 1ivb and 2qwk was only 32.07%. These results demonstrated that our system of classification and definition of distance were both reasonable.

The present clustering method successfully carried out the functional classification of proteins; in particular, proteins with low homology, as well as non-homologous proteins could be classified. One drawback of our classification method was observed in comparison with conventional sequence-based analysis. When a sequence-based classification system is applied, there is a correlation between the distance between two sequences and evolution time. In contrast, with our method, the distance between proteins depends on the docking score, and thus the biological significance of the distance remains unclear.

Fig. 1b shows some of the partial compound-classification results from the analysis of 31 pockets versus 132 compounds out of 132 pockets versus 132 compounds. The definition of distance was given by Eq. (4). The compounds

shown are the ligands of the protein–ligand complexes shown in Fig. 1a. The receptors are arbitrarily classified into three clusters for convenience. There was a correlation between the results in Fig. 1a and those in Fig. 1b; most of the ligands in cluster A in Fig. 1a were classified in cluster a in Fig. 1b, except for 2ctc and 1lna. For clusters b and c, there was no correlation between the compound cluster and the receptor cluster. Cluster b consisted of five ligands of cluster B, five ligands of cluster C, and three ligands of cluster D. Cluster c consisted of two ligands of cluster B, three ligands of cluster C, and one ligand of cluster D. The ligands carry some structural feature of the original docking pocket, for example, a pharmacophore. The correlation between the cluster of proteins and that of the ligands suggests that the ligands could be classified based on the kind of pharmacophore. This feature is different from the conventional compound-classification methods that are based on only the information of the compounds themselves.

There was a good correlation between the cluster and the number of atoms in the compound (see Fig. 2). Cluster a consisted primarily of small compounds. Typically, the number of atoms in the compounds in cluster a was 20–25, and most of these compounds were sugars. Cluster b consisted primarily of medium-sized compounds. The number of atoms in the compounds in cluster b tended to be 60–90, and there was a wide variety in compound structure. Cluster c consisted primarily of large compounds. The number of atoms in these compounds tended to be 100–110, and most of these compounds were peptide-mimetic compounds. Protein pockets recognize the volume of the ligand. A protein pocket that binds a small ligand cannot bind a compound that is larger than the pocket volume. Because our compound clustering method is based on the protein–ligand docking, this result is reasonable. This trend also agrees with the previous observation that the binding free energy increases with the number of non-hydrogen atoms. The larger ligand shows the stronger affinity rather than the smaller ligand, and these larger or smaller ligands with similar affinity form a cluster by our classification method [39].

The clustering of compounds was not as clear as that of receptors. Fig. 2 shows that the clustering of compounds strongly depends on the size of the compound. The use of our compound classification did not show a distinct advantage over a conventional similarity search, probably because the variety of the compounds in the data set was very large with almost no series of similar compounds. Our database consists primarily of proteases and peptidases. In general, even if there is a variety of putative receptor pockets, most receptors can bind peptide-like compounds, and certain common features will tend to exist between compounds. It is possible that this type of bias will prevent an unequivocal clustering of compounds.

This clustering method will also work when the affinity matrix calculated by the other docking software (DOCK, FlexX, GOLD, and etc.) is used. The detail of the clustering will change depending on the docking software, since the success rate of docking will be different depending on the docking software and the type of proteins.

### 4.2. Screening results

Fig. 3 shows database enrichments with the raw docking score, the MASC score, and the RS method applied to 132 receptors versus 132 compounds. For the RS method, the acceptance ratio $R_a$ and $\lambda$ were optimized to maximize the database enrichment; here, $R_a$ was set to 90%, and $\lambda$ was set to $-0.01$. In this study using the X-ray crystal structures of the protein–ligand complexes, the ligand in the complex structure was regarded as a hit compound for each protein. Usually, the in silico screening method is used to select only

Fig. 3. Database enrichment results of a total of 132 receptors vs. 132 compounds. The dashed line, empty squares, empty circles, and filled circles represent the results obtained with a uniform sampling, the raw docking score, the MASC score, and the RS method, respectively. The numbers of compounds and hits were scaled to %.
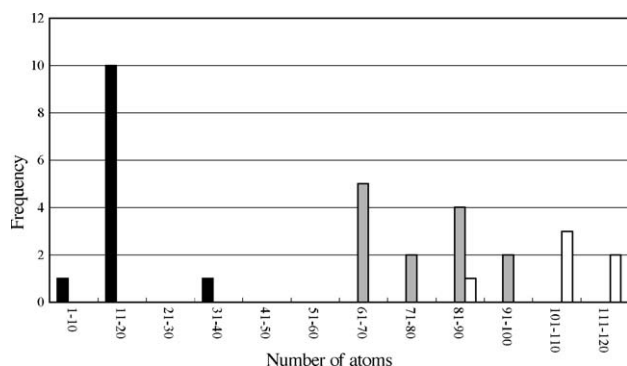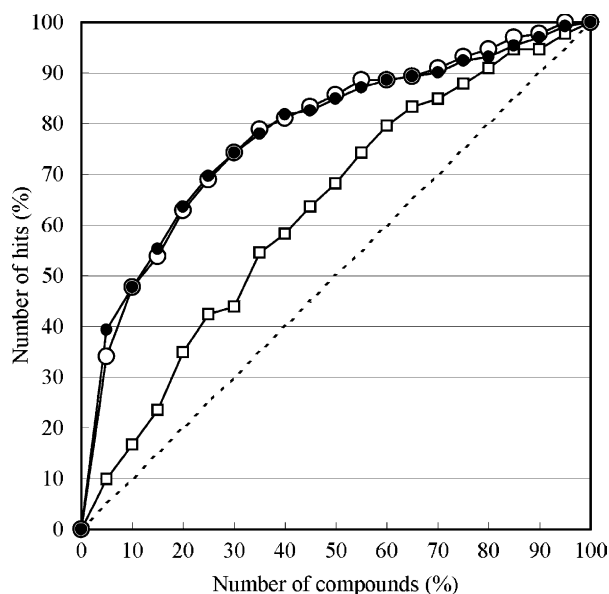
Fig. 2. Distribution of the number of atoms in the compounds. Black bars, grey bars, and white bars represent the distribution of the number of atoms in clusters a, b, and c, respectively.

a small number of compounds from a large number of compounds in a database, such that the number of hits obtained among the first small number of compounds of the database is crucial. Thus, in this study we have addressed only the number of hits found among the first 5 and 10% of the entries in the database. As shown in Fig. 3, the raw docking score yields a slight enrichment, with 9.8 and 16.6% of the ligands found among the first 5 and 10% of the database, respectively. The MASC score yielded a 6.8–4.8-fold enrichment, with 34.0 and 47.7% of the ligands found among the first 5 and 10% of the database, respectively. The RS method gives good enrichment, as was the case with the MASC scoring method. The RS method yielded a 7.9–4.8-fold enrichment, with 39.4 and 47.7% of the ligands found among the first 5 and 10% of the database, respectively.

Although our rough score function for docking was not optimized to reproduce the binding free energy, the score function worked well, as described above. The conventional screening methods with the raw scores compare the docking scores among different compounds that are composed of individual functional groups. Any score function depends on those functional groups so strongly that the accuracy of the docking scores would not be very high, causing the database enrichment to worsen. In contrast, multiple target methods, such as the MASC scoring method or the RS method, compare the interactions between one compound and many proteins. Namely, these methods attempt to determine which protein shows the highest affinity with the compound in question. Thus, the difference among those docking scores is expected to have high accuracy. The rate of the hit compounds is so rare in the current compound database, that the pair of a strongly interacting protein and the compound could be a hit pair. Consequently, the database enrichment in the multiple target method could become much better than that of the docking method with the raw scores.

The enrichment discussed above did not reach the theoretical upper limit, and could therefore still be improved. The docking program reconstructed 50.8 and 59.8% of the complexes with RMSD values of <2, and 3 Å, respectively. The enrichment was limited by the docking accuracy, and thus a more accurate docking program would yield higher levels of enrichment. Moreover, the docking score required to evaluate binding affinity could be improved. For example, our score did not include the term for the entropy change of rotational bonds in the ligand.

The advantages of the MASC score and the RS method over the raw-docking score depend greatly on the choice of receptors. Here, we constructed a new receptor set, i.e., a union of cluster 1 and cluster 7 of the total set of 132 receptors. The total number of receptors was 23, and all 132 compounds were used for the calculations. Fig. 4 shows the database enrichments obtained with the raw docking score, the MASC score, and the RS method. For the RS method, the acceptance ratio $R_a$ and $\lambda$ were optimized to maximize the database enrichment; $R_a$ was set to 50%, and $\lambda$ was set to −0.01. The raw docking score yielded a slight enrichment,
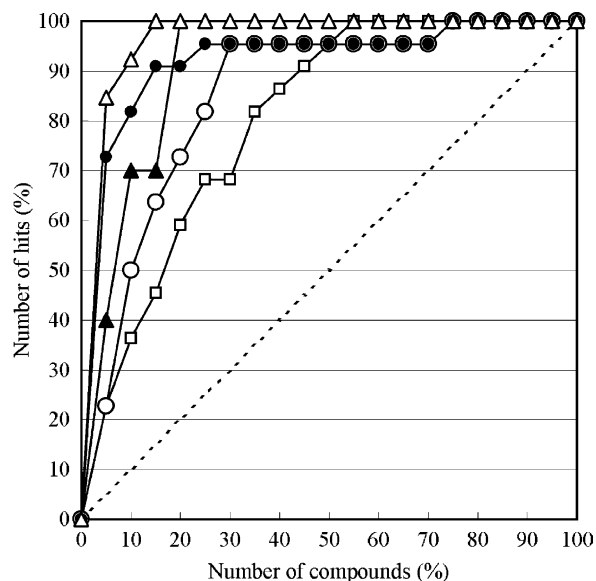


Fig. 4. Database enrichment results of 23 receptors, i.e., a union of cluster 1 and cluster 7, vs. 132 compounds. The dashed line, empty squares, empty circles, and filled circles represent the results obtained with a uniform sampling, the raw docking score, the MASC score, and the RS method, respectively. The filled and empty triangles represent the results obtained with the MASC scores of clusters 1 and 7, respectively. The numbers of compounds and hits were scaled to %.

with 22.7 and 36.3% of the ligands found among the first 5 and 10% of the database, respectively. The MASC score yielded a 4.5–5.0-fold enrichment, with 22.7 and 50.0% of the ligand found among the first 5 and 10% of the database, respectively. The results of the RS method were drastically improved over those obtained using the MASC score. The RS method yielded a 14.5–8.2-fold enrichment, with 72.7 and 81.8% of the ligands found among the first 5 and 10% of the database, respectively.

The increase of the number of proteins does not always improve the database enrichment. Rather, the choice of the protein data set is essential for the enrichment. Clusters 1, 7 and the union of these clusters consist of 10, 13 and 23 proteins. When the MASC scoring method was applied to clusters 1, 7 and their union, 40.0, 84.6 and 22.7% of the ligands were found among the first 5% of the database, respectively. The result of the MASC scoring method applied to the union of clusters 1 and 7 was worse than the results for either cluster 1 or 7, and it was close to the result of the method with the raw scores. We tried the same analysis for the combinations of clusters 1 and 2, that of 2 and 3, that of 3 and 4, that of 4 and 5, that of 5 and 6, and that of 6 and 7, and found the same phenomenon as that of 1 and 7. The combination of cluster 1 and 7 showed the most drastic difference between the results of the MASC scoring and the RS methods.

The MASC score is a deviation from a standard docking score that was estimated from the docking scores with many different proteins in the data set. If the protein data set includes several well-separated protein clusters, the data set

should have a biased character. This bias of the protein data set should reduce the accuracy of the standard docking scores. In contrast, if the protein data set includes only one protein cluster, all the estimated standard docking scores could carry the same errors, which are finally cancelled out. Thus the choice of proteins is important, and our clustering method is useful to make a valuable data set with critical evaluation of the above biases.

We reduced the number of compounds to 23 hit compounds from a total of 132 compounds. This alteration did not significantly change the trend of the database enrichments; moreover, the optimal $\lambda$ and $R_a$ did not change. These results suggest that the optimal $\lambda$ and $R_a$ are not dependent on the choice of the compounds, but rather on the choice of the receptor. Thus, when the method was applied experimentally instead of theoretically, we were able to determine the optimal $\lambda$ and $R_a$ using a limited number of known hit compounds; then, in cases involving a database of unknown compounds, the RS method could be applied using the same $\lambda$ and $R_a$.

On the other hand, the set of 132 receptors was divided into 23 clusters, and one representative receptor was selected from each cluster. Thus, we selected 23 representative receptors from among 132 receptors, and then carried out the same analysis. The database enrichment obtained by each method is shown in Fig. 5. In the case of the RS method, the acceptance ratio $R_a$ and $\lambda$ were optimized to maximize the database enrichment; $R_a$ was set to 90%, and $\lambda$ was set to −0.01. The raw docking score yielded a slight enrichment, with 21.7 and 26.1% of the ligands found among the first 5 and 10% of the database, respectively. The MASC score

yielded a 6.1–3.9-fold enrichment, with 30.4 and 39.1% of the ligands found among the first 5 and 10% of the database, respectively. The results obtained by applying the RS methods were more accurate than those obtained with the MASC score. The RS method yielded a 8.7–4.4-fold enrichment, with 43.5 and 43.5% of the ligands found among the first 5 and 10% of the database, respectively. The number of receptors was the same as that shown in Fig. 3, but the RS method was not found to have a significant advantage over the MASC score method.

According to our study, the enrichment obtained with the RS method was significantly greater than that obtained with random sampling and conventional screening with the raw docking score. The RS method produced similar or better results than the screening method with the MASC score, although the observed advantage of the RS method was strongly dependent on the choice of receptors. Choice of $\lambda$ affected the database enrichment to a level of approximately 2–3%. On the other hand, the inclusion of $R_a$ could drastically improve the level of database enrichment. These results suggest that the choice of the receptor was an important factor in achieving good enrichment; here, we were able to achieve good enrichment by using the $R_a$ for an arbitrarily selected receptor set.

Furthermore, the multiple-target screening method was shown to have an additional advantage over the conventional single-target screening method. Namely, since it is possible to determine the relative binding scores between receptors for each compound, the MASC score and the RS methods can be used to avoid compounds that bind to undesirable receptors, which in turn reduces the possibility of encountering side effects.
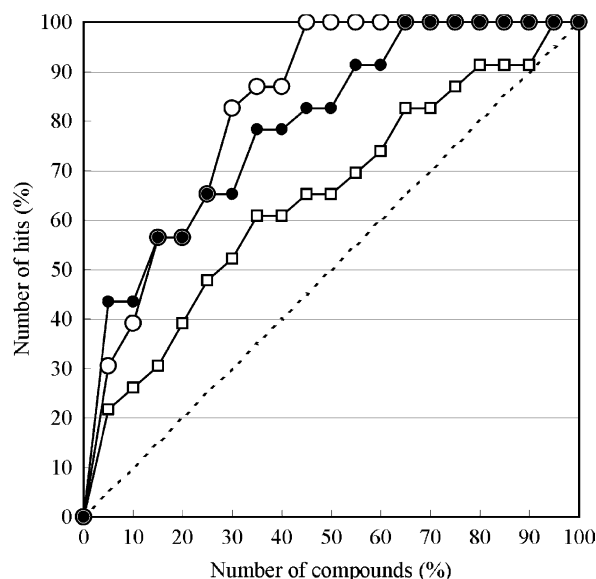


Fig. 5. Database enrichment results of 31 representative receptors vs. 132 compounds. The dashed line, empty squares, empty circles, and filled circles represent the results obtained with a uniform sampling, raw docking score, the MASC score, and the RS method, respectively. The numbers of compounds and hits were scaled to %.

## 5. Conclusions

The distances and similarities of receptor pockets and chemical compounds were determined based on a matrix of multi-receptor versus multi-ligand docking scores. The receptors were classified by a cluster analysis using the distance between receptors. Homologous proteins were classified in the same cluster, and proteins with low homology, as well as non-homologous proteins, that could bind to a similar ligand were also classified in the same cluster. These results revealed that both our system of classification and our definition of distance were reasonable. Moreover, the compounds analyzed in the present study were classified by cluster analysis. We were able to observe weak correlations between clusters of receptors and clusters of compounds.

Based on the line of reasoning presented here, we proposed a new in silico screening method, i.e., the receptor selection (RS) method, whereby compounds in a database are docked to multiple pockets. A receptor–ligand docking program was newly developed for this purpose.

In a test of the docking of 132 receptors versus 132 compounds, the RS method gave the best enrichment among

the screening methods using the raw docking score, the MASC score, and the RS method. The RS method yielded a 7.9–4.8-fold enrichment, with 39.4 and 47.7% of the ligands found among the first 5 and 10% of the database, respectively. The advantage of the RS method was found to depend greatly on the selection of receptors. In the multiple target method, the choice of proteins is essential, and the bias in the protein data set should be avoided. Our clustering method can thus be considered viable for constructing a suitable protein data set.

## Acknowledgements

## Appendix A

The following PDB identifier list of complexes was used: 1a28, 1a42, 1a4g, 1a4q, 1abe, 1abf, 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, 1ejn, 1epb, 1epo, 1ets, 1f0r, 1f0s, 1f3d, 1fen, 1fkg, 1fki, 1fl3, 1glp, 1hdc, 1hfc, 1hos, 1hpv, 1hsb, 1hsl, 1htf, 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, 1qpq, 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, 4est, 4lbd, 4phv, 5abp, 5cpp, 5er1, 6rnt, and 7tim. For 1abe, 1abf, 5abp, and 1htf, two receptor pockets were prepared, since these proteins both bind two ligands each.

## References

[1] C. Luong, A. Miller, J. Barnett, J. Chow, C. Ramesha, M.F. Browner, Flexibility of the NSAID binding site in the structure of human cyclooxygenase-2, Nat. Struct. Biol. 3 (1996) 927–933.

[2] P.N.P. Rao, M.J. Uddin, E.E. Knaus, Design, synthesis, and structure-activity relationship studies of 3,4,6-triphenylpyran-2-ones as selective cyclooxygenase-2 inhibitors, J. Med. Chem. 47 (2004) 3972–3990.

[3] X. Leval, J. Delarge, F. Somers, P. Tullio, Y. Henrotin, B. Pirotte, J.M. Dogne, Recent advances in inducible cyclooxygenase (COX-2) inhibition, Curr. Med. Chem. 7 (2000) 1041–1062.

[4] M. Miller, J. Schneider, B.K. Sathyanarayna, M.V. Toth, G.R. Marshall, L. Clawson, L. Selk, S.B. Kent, A. Wlodawer, Structure of complex of synthetic HIV-1 protease with a substrate-based inhibitor at 2.3 Å resolution, Science 246 (1989) 1149–1152.

[5] K. Kinoshita, H. Nakamura, Protein informatics towards function identification, Curr. Opin. Struct. Biol. 13 (2003) 396–400.

[6] K. Kinoshita, H. Nakamura, Identification of protein biochemical functions by similarity search using the molecular surface database eF-site, Protein Sci. 12 (2003) 1589–1595.

[7] A. Shulman-Peleg, R. Nussinov, H.J. Wolfson, Recognition of functional sites in protein structures, J. Mol. Biol. 339 (2004) 607–633.

[8] A. Schmitt, D. Kuhn, G. Klebe, A new method to detect related function among proteins independent of sequence and fold homology, J. Mol. Biol. 323 (2002) 387–406.

[9] S. Schmitt, M. Hendlich, G. Klebe, Development of a database for protein cavities and its usage for similarity searches in binding sites, in: H.D. Höltje, M. Sippl (Eds.), Rational Approaches to Drug Design, 2001, 135–141.

[10] C.T. Porter, G.J. Bartlett, J.M. Thornton, The catalytic site atlas: a resource of catalytic sites and residues identified in enzymes using structural data, Nucl. Acids Res. 32 (2004) D129–D133.

[11] L. Holm, C. Sander, Dali/FSSP classification of three-dimensional protein folds, Nucl. Acids Res. 25 (1997) 231–234.

[12] A. Heger, C.A. Wilton, A. Sivakumar, L. Holm, ADDA: a domain database with global coverage of the protein universe, Nucl. Acids Res. 33 (2005) D188–D191.

[13] R.A. George, R.V. Spriggs, J.M. Thornton, B. Al-Lazikani, M.B. Swindells, SCOPEC: a database of protein catalytic domains, Bioinformatics 20 (2004) i130–i136.

[14] I.G. Choi, J. Kwon, S.H. Kim, Local feature frequency profile: A method to measure structural similarity in proteins, Proc. Natl. Acad. Sci. U.S.A. 16 (2004) 3797–3802.

[15] K.A. Dennesiouk, M. Johnson, When fold is not important: A common structural framework for adenine and AMP binding in 12 unrelated protein families, Proteins: Struct., Funct. Genet. 38 (2000) 310–326.

[16] H.-J. Boehm, G. Schneider, R. Mannhold, H. Kubinyi, G. Folkers (Eds.), Protein—Ligand Interactions from Molecular Recognition to Drug Design—Methods and Principles in Medicinal Chemistry, Wiley-VCH, Weinheim, 2003, pp. 88–91.

[17] R.S. Pearlman, K.M. Smith, Metric validation and the receptor-relevant subspace concept, J. Chem. Inf. Comput. Sci. 39 (1999) 28–35.

[18] A.K. Ghose, V.N. Viswanadhan (Eds.), Combinatorial Library Design and Evaluation—Principle, Software Tools, and Applications in Drug Discovery, Marcel Dekker, New York, 2001, pp. 337–362.

[19] M.L. Verdonk, J.C. Cole, P. Watson, V. Gillet, P. Willett, SuperStar: improved knowledge-based interaction fields for protein binding sites, J. Mol. Biol. 307 (2001) 841–859.

[20] D.G. Truhlar, W.J. Howe, A.J. Hopfinger, J. Blaney, R.A. Dammkoehler (Eds.), Rational Drug Design, Springer-Verlag, New York, 1999, pp. 39–49.

[21] G.P.A. Vigers, J.P. Rizzi, Multiple active site corrections for docking and virtual screening, J. Med. Chem. 47 (2004) 80–89.

[22] I.D. Kuntz, J.M. Blaney, S.J. Oatley, R. Langridge, T.E. Ferrin, A Geometric approach to macromolecule–ligand interactions, J. Mol. Biol. 161 (1982) 269–288.

[23] M. Rarey, B. Kramer, T. Lengauer, G. Klebe, A fast flexible docking method using an incremental construction algorithm, J. Mol. Biol. 261 (1996) 470–489.

[24] G. Jones, P. Willet, R.C. Glen, A.R. Leach, R. Taylor, Development and validation of a genetic algorithm for flexible docking, J. Mol. Biol. 267 (1997) 727–748.

[25] N. Paul, D. Rognan, ConsDock: a new program for the consensus analysis of protein–ligand interactions, Proteins: Struct., Funct. Genet. 47 (2002) 521–533.

[26] C.A. Baxter, C.W. Murray, D.E. Clark, D.R. Westhead, M.D. Eldridge, Flexible docking using tabu search and an empirical estimate of binding affinity, Proteins: Struct., Funct. Genet. 33 (1998) 367–382.

[27] M.R. McGann, H.R. Almond, A. Nicholls, J.A. Grant, F.K. Brown, Gaussian Docking Functions, Biopolymers 68 (2003) 76–90.

[28] D.S. Goodsell, A.J. Olson, Automated docking of substrates to proteins by simulated annealing, proteins: structure, Funct. Genet. 8 (1990) 195–202.

[29] J.S. Taylor, R.M. Burnett, DARWIN: a program for docking flexible molecules, Proteins: Struct., Funct. Genet. 41 (2000) 173–191.

[30] R. Abagyan, M. Totrov, D. Kuznetsov, ICM: a new method for structure modeling and design: application to docking and structure prediction from the disordered native conformation, J. Comput. Chem. 15 (1994) 488–506.

[31] P.M. Colman, Structure-based drug design, Curr. Opin. Struct. Biol. 4 (1994) 868–874.

[32] Y. Lamdan, J. Schwartz, H. Wolfson, Affine invariant model-based object recognition, IEEE Trans. Robot. Automat. 6 (1990) 578–589.

[33] P.F.W. Stouten, C. Frommel, H. Nakamura, C. Sander, An effective solvation term based on atomic occupancies for use in protein simulations, Mol. Simul. 10 (1993) 97–120.

[34] J.W.M. Nissink, C. Murray, M. Hartshorn, M.L. Verdonk, J.C. Cole, R. Taylor, A new test set for validating predictions of protein–ligand interaction, Proteins: Struct., Funct. Genet. 49 (2002) 457–471.

[35] J. Wang, P. Cieplak, P.A. Kollman, How well does a restrained electrostatic potential (RESP) model perform in calculating conformational energies of organic and biological molecules? J. Comput. Chem. 21 (2000) 1049–1074.

[36] M.W. Schmidt, K.K. Baldridge, J.A. Boatz, S.T. Elbert, M.S. Gordon, J.H. Jensen, S. Koseki, N. Matsunaga, K.A. Nguyen, S. Su, T.L. Windus, M. Dupuis, J.A. Montgomery, The general atomic and molecular electronic structure system, J. Comput. Chem. 14 (1993) 1347–1363.

[37] M.J. Frisch, G.W. Trucks, H.B. Schlegel, G.E. Scuseria, M.A. Robb, J.R. Cheeseman, V.G. Zakrzewski, J.A. Montgomery, R.E. Stratmann Jr., J.C. Burant, S. Dapprich, J.M. Millam, A.D. Daniels, K.N. Kudin, M.C. Strain, O. Farkas, J. Tomasi, V. Barone, M. Cossi, R. Cammi, B. Mennucci, C. Pomelli, C. Adamo, S. Clifford, J. Ochterski, G.A. Petersson, P.Y. Ayala, Q. Cui, K. Morokuma, D.K. Malick, A.D. Rabuck, K. Raghavachari, J.B. Foresman, J. Cioslowski, J.V. Ortiz, A.G. Baboul, B.B. Stefanov, G. Liu, A. Liashenko P. Piskorz, I. Komaromi, R. Gomperts, R.L. Martin, D.J. Fox, T. Keith, M.A. Al-Laham, C.Y. Peng, A. Nanayakkara, C. Gonzalez, M. Challacombe, P.M.W. Gill, B. Johnson, W. Chen, M.W. Wong, J.L. Andres, C. Gonzalez, M. Head-Gordon, E.S. Replogle, J.A. Pople, Gaussian 98, Revision A.9; Gaussian, Inc., Pittsburgh PA, 1998.

[38] D.A. Case, T.A. Darden, T.E. Cheatham, III, C.L. Simmerling, J. Wang, R.E. Duke, R. Luo, K.M. Merz, B. Wang, D.A. Pearlman, M. Crowley, S. Brozell, V. Tsui, H. Gohlke, J. Mongan, V. Hornak, G. Cui, P. Beroza, C. Schafmeister, J.W. Caldwell, W.S. Ross, P.A. Kollman, AMBER 8, University of California, San Francisco, 2004.

[39] I.D. Kuntz, K. Chen, K.A. Sharp, P.A. Kollman, The maximal affinity of ligands, Proc. Nat. Acad. Sci. U.S.A. 96 (1999) 9997–10002.

[40] R.D. Page, TREEVIEW: an application to display phylogenetic trees on personal computers, Comput. Appl. Biosci. 12 (1996) 357–358.