# Prediction of binding constants of protein ligands: A fast method for the prioritization of hits obtained from de novo design or 3D database search programs

Hans-Joachim Böhm*

*BASF AG, Central Research, D-67056 Ludwigshafen, Germany*

## Summary

A dataset of 82 protein–ligand complexes of known 3D structure and binding constant $K_i$ was analysed to elucidate the important factors that determine the strength of protein–ligand interactions. The following parameters were investigated: the number and geometry of hydrogen bonds and ionic interactions between the protein and the ligand, the size of the lipophilic contact surface, the flexibility of the ligand, the electrostatic potential in the binding site, water molecules in the binding site, cavities along the protein–ligand interface and specific interactions between aromatic rings. Based on these parameters, a new empirical scoring function is presented that estimates the free energy of binding for a protein–ligand complex of known 3D structure. The function distinguishes between buried and solvent accessible hydrogen bonds. It tolerates deviations in the hydrogen bond geometry of up to 0.25 Å in the length and up to 30° in the hydrogen bond angle without penalizing the score. The new energy function reproduces the binding constants (ranging from $3.7 \times 10^{-2}$ M to $1 \times 10^{-14}$ M, corresponding to binding energies between $-8$ and $-80$ kJ/mol) of the dataset with a standard deviation of 7.3 kJ/mol corresponding to 1.3 orders of magnitude in binding affinity. The function can be evaluated very fast and is therefore also suitable for the application in a 3D database search or de novo ligand design program such as LUDI. The physical significance of the individual contributions is discussed.

## Introduction

Recently, several new computer programs for de novo ligand design [1–15] and for the docking of structures from 3D databases [16–25] have been described. These programs use a known 3D structure of the target protein to construct or retrieve a potentially very large number of possible protein ligands. An important prerequisite for a successful application of these new tools is the availability of scoring functions to prioritize the hits. In our opinion, this problem has not yet been solved satisfactorily.

The purpose of the present contribution is to describe our recent work aiming at the development of a fast scoring function that can prioritize a large list of up to several thousand structurally and functionally diverse ligands. The scoring function should be sufficiently accurate so that the testing of compounds can be limited to the top scoring structures. It should be applicable to a broad range of problems. Furthermore, it should be able to cope with small uncertainties in the 3D structure of the target protein.

There are basically two different situations where scoring functions play an important role in structure-based drug design. First, given that one has designed a possible ligand and has applied available docking algorithms to derive a reasonable 3D structure, the challenge is to predict the binding energy. In this case, the calculation time is not of primary importance. Therefore, methods like free energy pertubation theory and/or accurate calculations on the electrostatics [26–29] are possible approaches to calculate the binding

*Present address: Hoffmann-La Roche Ltd., Pharmaceuticals Division, Computational Chemistry, CH-4070 Basel, Switzerland.

energy. However, when applying methods for de novo design or for 3D database searching, the situation is different. Now, one is faced with the challenge to score and prioritize a large number of structurally diverse compounds. An important requirement is clearly that the scoring function be fast. It is this second problem that is addressed in the present paper.

Several methods have been proposed for scoring protein–ligand interactions [7, 8, 15, 30–44]. Some schemes use interaction energies obtained from a molecular mechanics force field [30–33], some base the prediction of binding affinities completely on the calculation of solvation energies [34], others use a combination of several terms including solvation energies [35]. Another approach to scoring is based on simple empirical functions [37, 42, 43]. The field has been reviewed by Ajay and Murcko [44].

Recently, a new simple scoring function was reported which was designed to predict binding constants for protein–ligand complexes of known three-dimensional structure [43]. This function (which in the present paper is referred to as SCORE1) takes into account hydrogen bonds, ionic interactions, the lipophilic contact surface and the number of rotatable bonds in the ligand. SCORE1 has been implemented in the de novo design program LUDI [14]. It was demonstrated that, when used together with LUDI, SCORE1 can successfully identify known high affinity ligands for the proteins trypsin and streptavidin out of a 3D database consisting of 30 000 small organic molecules [45]. The function was also able to predict correctly the very high affinity of a novel class of thrombin inhibitors [46]. It has also been used successfully in the recently developed program FlexX for docking of flexible ligands [47].

Since its development, we and others have applied SCORE1 to a large number of diverse protein–ligand complexes. This experience with SCORE1 has also led to the identification of protein–ligand complexes with deviations of more than three orders of magnitude (corresponding to errors in $\Delta G$ larger than 17 kJ/mol) between the predicted and the experimentally determined $K_i$ values. We have analyzed several protein–ligand complexes with high affinity ligands where our first scoring function underestimated the binding affinity. These complexes appear to be characterized by a number of characteristic features.

First, the visual inspection of several protein–ligand complexes with very tightly binding ligands revealed that these structures are characterized by a high steric complementarity. No cavities are present

along the protein–ligand interface. This phenomenon is not addressed in SCORE1. In the present work, the use of a cavity term is investigated that penalizes cavities along the protein–ligand interface.

Furthermore, the structures appear to exhibit a close-to-perfect electrostatic complementarity with no unpaired buried polar groups along the protein–ligand interface. This is also indicated by the fact that hydrogen acceptor groups of the ligand are located at positions where the electrostatic potential of the protein (calculated from a simple point charge model, see below) is positive, and donor groups are at positions with a negative electrostatic potential. This has prompted us to investigate the possibility of accounting for the electrostatic potential in protein–ligand interactions.

In the calibration of SCORE1 [43] the largest deviation between experimental and calculated binding constants was observed for the complex streptavidin–biotin. The poor performance of SCORE1 on this complex is thought to be due to an inadequate treatment of solvation effects. As pointed out previously [43], SCORE1 runs into problems if desolvation effects compensate for the direct protein–ligand interaction. As also discussed previously, another problem with SCORE1 is that the contributions from a buried and a solvent accessible hydrogen bond are scored identical. Therefore, we have tried to improve the treatment of solvation effects as will be discussed below.

In addition, our analysis revealed that several tight binding complexes including streptavidin–biotin achieve their hydrogen bonds with a comparatively small polar protein–ligand contact surface. This pointed to the possibility to use the ratio between the polar contact area and the number of the hydrogen bonds as a further parameter to distinguish between strong and weak hydrogen bonds.

A striking disagreement with experimental data was observed for the complex immunoglobulin–fluorescein where SCORE1 underestimated the binding affinity by 5 orders of magnitude in $K_i$. A visual inspection of this complex revealed an extensive network of interactions between aromatic rings of the protein and the ligand. The importance of aromatic interactions for ligand binding was pointed out previously [48–50]. SCORE1 does not contain a specific term for this type of interaction. Therefore, in the present work we have investigated the effect of adding such a term to the new scoring function.

The present approach is based on a geometrical analysis of 3D structures of protein–ligand complexes. We have sought to find a set of parameters that can be used to predict the binding affinity. We describe a new scoring function SCORE2, which is based on our previous function SCORE1, but attempts to address the problems described above. SCORE2 is based on a dataset consisting of 82 protein–ligand complexes.

## Methodology

The basic approach is the same as with our previous scoring function [43]. In comparison with SCORE1, we have augmented the scoring function by a number of additional terms which are highlighted by using bold characters. The following function is used:

$$\Delta G_{binding} = \Delta G_0 + \Delta G_{polar} + \Delta G_{apolar}$$
$$+ \Delta G_{solv} + \Delta G_{flexi} \tag{1}$$

$\Delta G_0$ is a contribution to the binding energy that does not directly depend on any specific interactions with the protein. It may be rationalized as a reduction of binding energy due to overall loss of translational and rotational entropy of the ligand. $\Delta G_{polar}$ and $\Delta G_{apolar}$ represent the polar and apolar interactions, $\Delta G_{solv}$ accounts for desolvation effects and $\Delta G_{flexi}$ treats the ligand flexibility. There is no change in the treatment of the ligand flexibility as compared to SCORE1:

$$\Delta G_{flexi} = \Delta G_{rot} NROT \tag{2}$$

$\Delta G_{rot}$ describes the loss of binding energy due to freezing of internal degrees of freedom in the ligand. The number of rotatable bonds NROT is taken as the number of acyclic $sp^3$-$sp^3$ and $sp^3$-$sp^2$ bonds. Rotations of terminal -$CH_3$ or -$NH_3$ groups are not taken into account.

*Polar interactions*
The polar contributions to $\Delta G_{binding}$ are described as follows.

$$\Delta G_{polar} = \Delta G_{hb} \Sigma_{h-bonds} f(\Delta R, \Delta\alpha)$$
$$\times \mathbf{f(N_{neighb})} \times \mathbf{f_{pcs}}$$
$$+ \Delta G_{ionic} \Sigma_{ionic\ int.} f(\Delta R, \Delta\alpha)$$
$$\times \mathbf{f(N_{neighb})} \times \mathbf{f_{pcs}} \tag{3}$$
$$+ \mathbf{\Delta G_{esrep} N_{repulsive\ contacts}}$$

$$f(\Delta R, \Delta\alpha) = f1(\Delta R)f2(\Delta\alpha)$$

$$f1(\Delta R) = \begin{cases} 1 \text{ if } \Delta R \leq TOL \\ 1 - (\Delta R - TOL)/0.4 \\ \qquad \text{if } \Delta R \leq 0.4 + TOL \\ 0 \text{ if } \Delta R > 0.4 + TOL \end{cases}$$

$$f2(\Delta\alpha) = \begin{cases} 1 \text{ if } \Delta\alpha < 30° \\ 1 - (\Delta\alpha - 30)/50 \text{ if } \Delta\alpha \leq 80° \\ 0 \text{ if } \Delta\alpha > 80° \end{cases}$$

$\Delta G_{hb}$ describes the contribution from an ideal hydrogen bond. $\Delta G_{hb}$ accounts for all hydrogen bonds with at least one partner being neutral. $\Delta G_{ionic}$ represents the contribution from an unperturbed ionic interaction. The same geometric dependency is assumed for uncharged and for charged interactions.

$f(\Delta R, \Delta\alpha)$ is a penalty function which accounts for large deviations of the hydrogen bond geometry from ideality. The functional form of $f(\Delta R, \Delta\alpha)$ is basically the same as in SCORE1 [43] with the exception that the tolerated deviation in the hydrogen bond length TOL is now 0.25 Å. $\Delta R$ is the deviation of the H··O/N hydrogen-bond length from the ideal value 1.9 Å. $\Delta\alpha$ is the deviation of the hydrogen bond angle $\angle_{N/O-H··O/N}$ from its idealized value of 180°. The function tolerates small deviations of up to 0.25 Å and 30° from the ideal geometry which are often due to small uncertainties in the X-ray structure. Similar to SCORE1, the present function does not include a term accounting for distortions in the hydrogen bond angle $\angle_{C=O··O/N}$. For interactions with metal atoms $f2(\Delta\alpha)$ is set to one and the ideal distance Metal··O/N is set to 2.1 Å.

When SCORE1 was applied to protein–ligand complexes optimized by a molecular mechanics force field calculation, it was observed that the predicted binding affinities were slightly overestimated. As a result of the force field optimization, the hydrogen bond geometries tend to move closer to their optimal values resulting in higher scores. We conclude that the allowed tolerances in SCORE1 are too small. We have therefore explored the use of larger tolerances to make the scoring function less sensitive to small deviations in the nonbonded contact geometries which are believed to result from experimental uncertainties without physical significance.

In a parameter study, we used the allowed deviation from ideality as an adjustable parameter in the scoring function. However, this fit gave even smaller allowed tolerances with a minimal standard deviation at TOL = 0.15 Å. This behaviour of the scoring function can be explained by a number of high resolution

structures with very strongly binding ligands present in the calibration dataset that contain hydrogen bond geometries close to ideality. The X-ray structures of the three most strongly binding ligands all contain several hydrogen bonds with geometries close to ideality. The goal of the present work was, however, to develop a scoring function that is also applicable to protein structures with resolutions higher than 2.5 Å. Therefore, it was decided to use a slighly larger tolerance at the expense of a slightly larger standard deviation.

We have also investigated a scoring function which further differentiates the polar interactions into six classes:

| # | Ligand | Protein |
|---|--------|---------|
| 1 | neutral | neutral |
| 2 | charged | neutral |
| 3 | neutral | charged |
| 4 | charged | charged |
| 5 | neutral | metal ion |
| 6 | charged | metal ion |

However, no improvement was obtained in the fit. For the interaction types #1, 2, 3 and 5, all $\Delta G$ values were within $\pm 20\%$ of the value averaged over all classes. Similarly, the difference between the contributions #4 and #6, as obtained from the fit, was less than 20%. In addition, we have also tested other schemes to further partition the polar interactions. The only scheme that actually improved the scoring function slightly, was to add a third term, which specifically accounts for interactions between a metal ion and a charged ligand atom. However, in view of the very small improvement of the fit, this scheme was not further investigated.

$f(N_{Neighb})$ is a new empirical function which serves to distinguish between convex and concave parts of the protein surface. The idea behind this function is to assign a higher weight to polar interactions in pockets as compared to those formed at the outer surface of the protein. We have adapted a simple scheme proposed by Sander [51]. For any given protein atom, $N_{Neighb}$ is simply taken as the number of non-hydrogen protein atoms that are closer than 5 Å. At present we use the following function:

$$f(N_{Neighb}) = (N_{Neighb}/N_{Neighb,0})^{\alpha}; \; \alpha = 1/2 \quad (4)$$

$f(N_{Neighb})$ is a property of each protein atom. A parameter study revealed an optimum at $\alpha = 1/2$. This value was then used in some of the scoring functions presented below. $N_{Neighb,0}$ was set to 25 for normal-

ization purposes. This value roughly corresponds to the average number of neighbours found in a subset of the calibration dataset. In the present dataset, $N_{Neighb}$ varies between 7 and 33. Therefore, $f(N_{Neighb})$ varies between 0.53 and 1.15.

In order to account for the polar contact surface area, we use the polar contact surface per hydrogen bond $f_{pcs}$ as a second factor to differentiate between strong and weak hydrogen bonds. Several different analytical forms for $f_{pcs}$ were investigated. Finally, it was decided to use a simple stepwise function:

$$f_{pcs} = \beta; \; A_{polar}/N_{HB} < 10 \; \text{Å}^2$$
$$f_{pcs} = 1; \; A_{polar}/N_{HB} > 10 \; \text{Å}^2 \quad (5)$$

$A_{polar}$ is the size of the polar protein–ligand contact surface and $N_{HB}$ is the number of hydrogen bonds. If $A_{polar}/N_{HB}$ is smaller than 10 Å$^2$ per hydrogen bond then $\Delta G_{hb}$ is multiplied by a factor $\beta$. A parameter study on $\beta$ revealed an optimum at 1.2. This value was then used in some of the scoring functions described below. $f_{pcs}$ is a property of a protein–ligand complex.

In addition, we have investigated the use of the electrostatic potential $V_{es}$ as an additional factor modulating the contribution of the polar interactions to the binding affinity. A simple point charge model was used to estimate the electrostatic potential of the protein in the binding site. The point charge model was adapted from the GROMOS force field [52]. $V_{es}$ was calculated on a cubic grid with 1 Å grid spacing. The value at the position of the polar ligand atom was then determined by linear interpolation from the surrounding eight points of the grid. Note that, in contrast to traditional force fields, the point charge model is not used to calculate an electrostatic interaction energy. It merely serves to put a higher weight on those hydrogen bonds and ionic interactions where $V_{es}$ strongly deviates from zero. $V_{es}$ is calculated using an effective distance dependent dielectric as also suggested previously by others [35]. A cutoff of 10 Å was used in the calculations.

$$V_{es} = \Sigma_{protein \; atoms} q_i/\epsilon r, \; \epsilon = 4r \quad (6)$$

All attempts to improve the scoring function by using a point-charge model for the electrostatic interactions were unsuccessful. Even with the simplest model, just taking the sign of $V_{es}$ (if $V_{es}$ is negative at the position of a ligand donor group or positive at the position of a ligand acceptor group, the interaction is assumed to be stronger than other interactions), a poor fit to the data was obtained. An analysis of the electrostatic potential at the position of the polar ligand atoms

revealed that more than 90% of the polar atoms are at positions where $V_{es}$ has the correct sign. However, a much higher rate of mismatches was observed for proteins with a large overall charge. We conclude that a simple point charge model, that neglects polarisation effects and charge compensation by counter ions, is not useful in a simple scoring function if the dataset contains several different proteins.

It has recently been shown that the relative binding affinities of some host-guest complexes can be understood by considering secondary electrostatic interactions [53]. These are interactions between polar groups that do not form a direct hydrogen bond but are relatively close to each other. We reasoned that a similar scheme could also be useful in a scoring function for protein–ligand complexes. Therefore, the additional parameter $\Delta G_{esrep}$ was included. In the present implementation, only repulsive interactions $\Delta G_{esrep}$ are taken into account. If two donor groups approach each other with a H-H distance of less than 3 Å or two acceptor groups have an O/N-O/N distance smaller than 4 Å, they are taken into account as one repulsive interaction. The contribution from repulsive electrostatic interactions is assumed to be proportional to the number of those contacts $N_{esrep}$.

*Apolar interactions*

The apolar interactions are described by two terms:

$$\Delta G_{apolar} = \Delta G_{lipo}|A_{lipo}| \\ + \Delta G_{aro} \Sigma_{aro-int.} f(R, \Theta) \qquad (7)$$

$\Delta G_{lipo}$ represents the contribution from lipophilic interactions. $A_{lipo}$ is calculated as described previously [43]. The lipophilic interaction is assumed to be proportional to the protein–ligand contact surface area. We now use a revised set of atomic radii. The following atomic radii are used for both the ligand and the protein: C: 1.8 Å, H: 1.2 Å, O,N: 1.6 Å, F: 1.5 Å, Cl: 1.8 Å, Br: 1.9 Å, S: 2.0 Å. All grid points inside an atomic sphere with the radius given above +0.4 Å are marked as occupied by the ligand. The same increased radii are used to mark the cubes that are adjacent to the ligand and overlap with the protein. The increment of 0.4 Å was derived in a trial-and-error procedure using a small subset of the protein–ligand complexes given in Table 1. The final calculated surface area is obtained by multiplying the number of surface cubes with a calibration factor (0.88) which was derived by comparing the present results with data obtained from Connolly's MS program [54]. It should be noted that the present approach to calculate the lipophilic surface

is also error-tolerant. Slight overlaps or small gaps between the protein and the ligand are ignored and do not affect the size of the calculated surface area.

In addition, the use of two different types of lipophilic contacts was investigated. $A_{lipo}$ was separated into two terms $A_{lipo1}$ and $A_{lipo2}$. $A_{lipo1}$ is the lipophilic contact surface where both the ligand and the protein are lipophilic. $A_{lipo2}$ is the contact surface where the ligand is lipophilic and the protein is polar or vice versa. The concept of interaction sites [14] was used to differentiate between the lipophilic and the polar part of the protein surface. All cubes that are at suitable positions for forming a polar interaction with the protein are treated as polar. The remaining ones are treated as lipophilic. Unfortunately, no improvement of the scoring function was obtained and the concept was therefore not further pursued.

$\Delta G_{aro}$ is a new term which describes interactions between aromatic rings. The importance of this interaction for protein structures, protein–ligand complexes and host-guest complexes was pointed out previously by others [48–50]. A detailed statistical analysis of aromatic contacts in high resolution protein structures revealed only a slight preference for the electrostatically favorable interaction geometries [50] (aromatic interactions are therefore not purely apolar). We have decided to ignore the angular dependence of this interaction at present and use a simple distance cutoff. If the distance between any heavy atoms of the rings is smaller than 4.5 Å, the rings are assumed to interact.

*Desolvation effects*

In view of the apparent importance of desolvation effects in protein–ligand interactions, several different schemes were investigated to account for this effect. First, we used a desolvation model with additive surface contribution. However, this model did not yield a reduced standard deviation of the prediction. Furthermore, the parameters obtained from the fit appeared physically unrealistic. For example, large negative contributions to $\Delta G_{binding}$ were obtained for the removal of polar parts of the surface from solvent, whereas measured solvation energies would suggest the opposite to be the case. If, on the other hand, the parameters for the desolvation model were fixed at values taken from the literature, a poor fit to the binding data was obtained.

We then tried to use two new terms, which represent the desolvation of polar ligand atoms. We simply used the number of polar and charged ligand atoms

*Table 1.* Protein–ligand complexes used for the calibration of the free energy function. The $-\log K_i$ predicted values refer to results obtained with function #7

| Protein–ligand complex | PDB entry | $-\log K_i$ Pred. | $-\log K_i$ Exp. | Ref. |
|---|---|---|---|---|
| Adenosinedeaminase – deazaadenosine | 1ADD | 6.61 | 6.74 | 59 |
| Carbonic anhydrase – methazolamide | 1BZM | 4.87 | 6.03 | 60 |
| Carboxypeptidase – benzylsuccinate | 1CBX | 7.50 | 6.30 | 61 |
| Carboxypeptidase – sulfodiimide | 1CPS | 6.56 | 6.66 | 62 |
| Cytidine deaminase – 4 dehydrozebularine | 1CTT | 3.69 | 4.52 | 63 |
| Elastase – TFA-Lys-Pro-*p*-isopropylanilide | 1ELA | 7.15 | 6.35 | 64 |
| Elastase – TFA-Lys-Phe-*p*-isopropylanilide | 1ELC | 5.35 | 7.15 | 64 |
| FKFB – FK506 | 1FKF | 7.41 | 9.70 | 65 |
| HIV protease – VX478 | 1HPV | 8.22 | 9.22 | 66 |
| HIV protease – XK263 | 1HVR | 9.19 | 9.51 | 67 |
| Lysozyme (L99A mutant) – benzofuran | 1L82 | 4.34 | 3.95 | 68 |
| Lysozyme (L99A mutant) – benzene | 1L83 | 3.69 | 3.75 | 68 |
| Lysozyme (L99A mutant) – phenylbutane | 1L86 | 6.20 | 4.85 | 68 |
| Lysozyme (L99A mutant) – *p*-xylene | 1L87 | 5.09 | 3.37 | 68 |
| Lactate dehydrogenase – oxamate | 1LDM | 5.14 | 5.40 | 69 |
| Myoglobin – imidazole | 1MBI | 3.60 | 1.88 | 70 |
| Cytochrome P450-2-phenyl-imidazole | 1PHE | 4.18 | 5.70 | 71 |
| Cytochrome P450-4-phenyl-imidazole | 1PHF | 5.50 | 4.40 | 71 |
| Cytochrome P450-metyrapone | 1PHG | 7.01 | 8.66 | 71 |
| Trypsin – thrombstop | 1PPC | 7.77 | 6.46 | 72 |
| Trypsin – 3-TAPAP | 1PPH | 7.59 | 6.22 | 72 |
| Pepsin – pepstatine | 1PSO | 8.64 | 10.34 | 73 |
| Retinol binding protein – retinol | 1RBP | 5.78 | 6.72 | 74 |
| Renin – CGP38560 | 1RNE | 9.29 | 9.40 | 75 |
| Sulfate binding protein – sulfate | 1SBP | 4.26 | 6.92 | 76 |
| Streptavidin – HABA | 1SRE | 6.13 | 4.00 | 77 |
| Streptavidin – biotin | 1STP | 10.85 | 13.40 | 77 |
| Thermolysin – phosphoramidon | 1TLP | 7.15 | 7.55 | 78 |
| Thermolysin – CLT | 1TMN | 10.30 | 7.30 | 78 |
| Trypsin – phenylpropylamine | 1TNK | 2.82 | 1.49 | 31 |
| Cytochrome P450 – camphore | 2CPP | 4.78 | 6.07 | 79 |
| Carboxypeptidase A-phenyllactate | 2CTC | 5.80 | 3.89 | 80 |
| Endothiapepsin – H256 | 2ER6 | 7.48 | 7.22 | 81 |
| Galactose binding protein – galactose | 2GBP | 6.08 | 7.60 | 82 |
| Glycogen phosphorylase – glucose | 2GPB | 3.64 | 2.77 | 83 |
| Fatty acid binding protein – $C_{15}COOH$ | 2IFB | 5.67 | 5.43 | 84 |
| Hermolysin – PLN | 2TMN | 4.16 | 4.67 | 78 |
| Thymidylatesynthase – CB3717 | 2TSC | 7.02 | 8.52 | 85 |
| Thymidylatesynthase – Cmpd. 3 | 2TSC[f] | 6.54 | 5.40 | 85 |
| Xyloseisomerase – xylitol | 2XIS | 5.13 | 5.82 | 73 |
| TIM – phosphoglycolic acid | 2YPI | 3.93 | 4.82 | 73 |
| Carboxypeptidase A – GT | 3CPA | 3.77 | 3.88 | 86 |
| DHFR – methotrexate | 3DFR | 9.50 | 10.3 | 87 |
| Trypsin – benzamidine | 3PTB | 4.81 | 4.74 | 88 |
| Trypsin – butylamine | 3PTB[a] | 3.88 | 2.82 | 88 |
| Trypsin – benzylamine | 3PTB[a] | 4.00 | 3.42 | 88 |
| Trypsin – phenylguanidine | 3PTB[a] | 4.56 | 4.14 | 88 |
| Trypsin – BPTI – IleVal | 3TPI | 6.02 | 4.30 | 89 |

*Table 1.* Continued

| Protein–ligand complex | PDB entry | −log $K_i$ Pred. | −log $K_i$ Exp. | Ref. |
|---|---|---|---|---|
| Chymotrypsin – benzene | 4CHA[a] | 2.20 | 1.60 | 90 |
| Chymotrypsin – phenole | 4CHA[a] | 2.69 | 2.19 | 90 |
| Chymotrypsin – indole | 4CHA[a] | 3.06 | 3.10 | 90 |
| Chymotrypsin – benzoquinoline | 4CHA[a] | 3.76 | 4.20 | 90 |
| Concanavalin – α-me-mannosid | 4CNA | 3.62 | 2.00 | 91 |
| DHFR – methotrexate | 4DFR | 9.24 | 9.70 | 87 |
| DHFR – 2,4-diaminopteridine | 4DFR[f] | 4.57 | 6.00 | 92 |
| Endothiapepsin – ac-pepstatin | 4ER2[e] | 6.89 | 8.04 | 73 |
| Endothiapepsin – H142 | 4ER4 | 8.09 | 6.79 | 93 |
| Gluthathione reductase – rGSSGr | 4GR1 | 0.99 | 1.70 | 94 |
| HIV protease – MVT101 | 4HVP | 7.35 | 6.15 | 95 |
| HIV protease – L700, 417 | 4PHV | 9.87 | 9.15[d] | 96 |
| Thermolysin – Leu-NHOH | 4TLN | 3.88 | 3.72 | 78 |
| Thermolysin – ZF$^P$LA | 4TMN | 10.26 | 10.19 | 78 |
| Tyr transferase – Tyr | 4TS1 | 4.27 | 5.60 | 97 |
| Cytochrome P450 – adamantone | 5CPP | 5.05 | 5.88 | 79 |
| Triose phosphate isomerase – $SO_4$ | 5TIM | 1.38 | 2.30 | 98 |
| Thermolysin – HONH-BAGN | 5TLN | 6.17 | 6.37 | 78 |
| Thermolysin – ZG$^P$LL | 5TMN | 8.17 | 8.04 | 78 |
| Thermolysin – thiorphan | 5TMN[c] | 5.18 | 5.64 | 99 |
| Aconitase – tricarballylate | 6ACN | 2.53 | 3.00 | 100 |
| Carboxypeptidase A – ZAA$^P$(O)F | 6CPA | 9.24 | 11.52 | 101 |
| Ribonuclease – uridinevanadate | 6RSA | 5.54 | 5.00 | 102 |
| Carboxypeptidase – ZFV$^P$(O)F | 7CPA | 11.92 | 14.0 | 103 |
| Cytochrome P450 – norcamphore | 7CPP | 4.28 | 3.80 | 79 |
| Asp aminotransferase – 2-pyridoxal-5$PO_4$ | 9AAT | 11.07 | 8.22 | 104 |
| Triose phosphate isomerase – GAPDH | _[j] | 4.82 | 4.70 | 105 |
| Thrombin – NAPAP | _[b] | 9.18 | 8.52[g] | 106 |
| Thrombin – LU57348 | _[i] | 10.12 | 10.30[h] | 46 |
| Thrombin – TAPAP | _[b] | 6.61 | 6.19 | 107 |
| Thrombin – amidinopiperidine | 1DWB[a] | 4.71 | 3.82 | 108 |
| Glycogen phosphorylase – cmpd #3 | _[i] | 2.51 | 1.43 | 109 |
| Glycogen phosphorylase – cmpd #10 | _[i] | 5.05 | 4.09 | 109 |
| Glycogen phosphorylase – cmpd #11 | _[i] | 3.89 | 5.52 | 109 |

[a] The ligand was positioned into the protein binding site using the program LUDI. See Ref. 43 for a discussion.

[b] The protein structure determined by Brandstetter et al. [110] was used.

[c] The protein coordinates were taken from 5TMN and the ligand coordinates were taken from Ref. 99. The S··Zn interaction present in this complex was counted as an ionic interaction.

[d] $IC_{50}$ value.

[e] The ligand was taken from the X-ray structure and the acetyl group was appended.

[f] The ligand was taken from the X-ray structure and part of the ligand was removed in order to obtain the specified structure.

[g] The reported $K_i$ of 6 nM was obtained for the racemate. We used $K_i = 3$ nM as only one enantiomer binds to the protein.

[h] The protein structure determined by Mack et al. [46] was used. The experimentally determined $IC_{50}$ is 60 times lower than the $IC_{50}$ value of NAPAP. It was assumed that the $K_i$ is also 60 times lower than that of NAPAP.

[i] The coordinates were kindly provided by Dres. L. Johnson and K. Watson, University of Oxford.

[j] The coordinates were kindly provided by Dres. W. Hol and C. Verlinde, University of Washington.

which are removed from solvent upon ligand binding as a measure for desolvation. $\Delta G_{solv,polar}$ and $\Delta G_{solv,charged}$ are the desolvation energies that are required for desolvation of the corresponding atoms. These terms should be adverse to ligand binding. Again, the inclusion of these two terms in the scoring function resulted in unrealistic values for the parameters and no improvement of the fit was observed.

Finally, we investigated the use of explicit waters in the binding site to account at least for some contribution to desolvation. The protein binding site was filled with water molecules. The positions of the waters were optimized using 100 steps of steepest descent minimization, followed by 1000 steps of conjugate gradients minimization. Then, a 1 ps molecular dynamics simulation was carried out followed by an additional minimization. The protein was kept rigid. The binding site was then visually inspected. If regions in the binding site with no water molecules were detected, the corresponding holes were filled with additional waters and the complete procedure was repeated. All calculations were performed with the programs INSIGHT and DISCOVER [55] using the CVFF force field [56]. Finally, only those water molecules were taken into account, that coincide with the occupied volume of the bound ligand. It was decided to distinguish between water molecules that form one or more hydrogen bonds with the protein and those that do not form a hydrogen bond. The inclusion of these two parameters, $\Delta G_{hb\ water}$ and $\Delta G_{lipo\ water}$, resulted in models with a significantly reduced standard deviation. However, the contribution from $\Delta G_{hb\ water}$ was very small. This term was therefore neglected and only $\Delta G_{lipowater}$ was used as a parameter in SCORE2.

$$\Delta G_{solv} = \Delta G_{lipo\ water} \Sigma \text{unbound water molecules}$$
(8)

Obviously, deriving the number of unbound water molecules for a protein cavity is time consuming, although it needs to be done only once.

*Investigation of a cavity term*
Guided by the visual inspection of several complexes with tightly binding ligands, it was decided to investigate the use of a cavity term. $\Delta G_{cavity}$ is a new term that penalizes cavities along the protein–ligand interface. A flood-fill algorithm as also used in the program VOIDOO [57] was used to calculate cavities along the protein/ligand interface. It was decided to penalize only lipophilic cavities. A new term $\Delta G_{cavity}\ V_{cavity}$ was added to the scoring function. Disappointingly,

a fit to the experimental data revealed a very small contribution for $\Delta G_{cavity}$ with no improvement for the standard deviation. An analysis of the dataset revealed that for the majority of the protein–ligand complexes, $V_{cavity}$ is smaller than 30 Å$^3$ which is comparable to the volume of a water molecule. Therefore, most of the cavities are smaller than a water molecule. In other words, it appears that, with some exceptions, there are basically no lipophilic cavities along the protein–ligand interface. However, we cannot exclude the possibility that the results from our calculation are just not accurate enough to allow for a meaningful calibration of the term $\Delta G_{cavity}$. It was decided to refrain from using this term in the scoring function.

*Selection of the calibration dataset*
The list of protein–ligand complexes used in the present study to calibrate the scoring function is given in Table 1. The coordinates of the protein–ligand complex were taken from the Brookhaven protein databank (PDB) [58] if not specified otherwise. We have selected 82 structures with known $K_i$ values [31, 46, 59–109] as basis for the calibration of the energy function. In the present work, we have only included structures with a reported resolution better than 2.7 Å. The structures were taken as stored in the PDB. Water molecules were excluded. Hydrogens were added using the graphics program INSIGHT [55]. For the side chain hydroxyl groups of Ser, Thr and Tyr the hydrogen atoms were positioned according to the hydrogen bond pattern observed in the structure. In all other cases (when the position of the hydrogen could not be deduced from the surrounding atoms), the hydrogen atoms were positioned in the trans orientation with respect to H-O-C-C. No energy minimization was carried out on any of the structures used in the present study. The amino acids Asp, Glu, Lys and Arg were assumed to be charged if not stated otherwise by the authors who determined the structure. Histidine side chains were protonated as indicated in the original publications. Accordingly, in some cases His was assumed to be charged. In the ligands the following groups were assumed to be charged: phosphate, phosphonamide, phosphonate, carboxylate, guanidinium, amidinium and amine. As discussed previously [43], the ureido group of biotin complexed with streptavidin was treated as a zwitter ion. The experimental binding data were taken from the literature without any further modifications. We did not check the experimental data for differences in temperature or salt concentrations during measurement. In the selection

of the structures, we have attempted to cover a broad spectrum of different types of protein–ligand complexes. The experimentally observed $K_i$ values range from 37 mM to 10 fM and thus cover more than 12 orders of magnitude.

The present dataset contains the 54 structures from our previous work with the exception of the structures 1DWB, 1ULB, 2R04, 2PHH and 9HVP which were excluded due to insufficient resolution (except for 2PHH, SCORE1 reproduced the $K_i$ of these structures very well with an error of less than one log unit). Thirty-three new protein–ligand complexes have been added. The present dataset is now significantly larger than the previous one and covers a broader range of structures. However, there is still a bias towards serine proteases, aspartic proteases and metallo proteases due to their dominant presence in the PDB. Furthermore, a large part of the structures contains ligands that are peptidic in nature. The number of ligands without amide bonds and with heterocyclic moieties is still very small.

The experimental $K_i$ values were converted into free energies of binding $\Delta G$ using $\Delta G = -RT \ln K_i$ (T = 298 K). A least squares fit was then performed to obtain the adjustable parameters in the free energy function.

## Results

The aim of the present work was to develop an efficient scoring function for the estimation of binding affinities of protein–ligand complexes. We have investigated several variants of the scoring functions given above. The results are summarized in Table 2. Function #1 is identical to SCORE1 and was taken directly from our previous work [43]. It contains four adjustable parameters, $\Delta G_{hb}$, $\Delta G_{ionic}$, $\Delta G_{lipo}$, $\Delta G_{rot}$, and the constant term $\Delta G_0$. SCORE1 reproduces the current data with a standard deviation of 9.5 kJ/mol corresponding to 1.6 orders of magnitude in $K_i$. Function #2 has the same analytical form as SCORE1, with the adjustable parameters recalibrated to the new extended dataset. The new calibration reduces the standard deviation to 8.8 kJ/mol or 1.5 orders of magnitude in $K_i$. The recalibration does not yield a significant improvement, indicating that SCORE1 is also applicable to a broad range of structures. Function #3 has also the same analytical form as #1 and #2 with the exception that the new atomic radii described above are used to
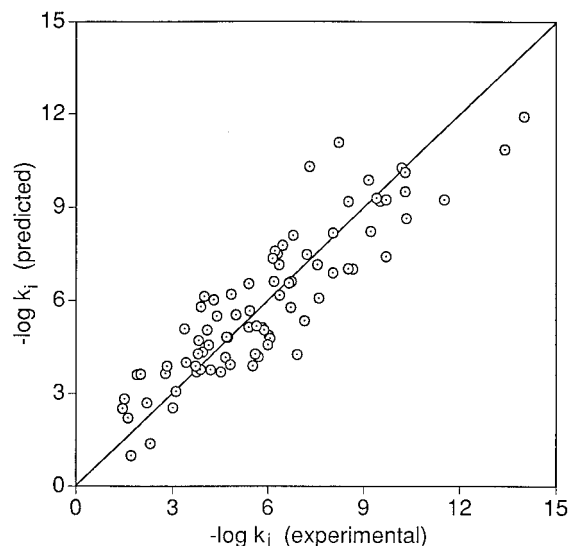


*Figure 1.* Plot of the calculated binding constants $K_i$ from function #7 versus experimentally observed values for the 82 protein–ligand complexes of the calibration dataset.

calculate the lipophilic contact surface. The standard deviation is slightly reduced to 8.6 kJ/mol.

The effect of the inclusion of the term $\Delta G_{aro}$ is seen in function #4, which contains five adjustable parameters. The standard deviation is 8.1 kJ/mol. The further inclusion of the parameter $\Delta G_{lipowater}$ yields function #5 with six adjustable parameters and a standard deviation of 7.5 kJ/mol. Finally, the inclusion of $\Delta G_{esrep}$ improves the fit only marginally by 0.1 kJ/mol. However, because the actual value for $\Delta G_{esrep}$ obtained from the fit is in line with literature values for the contribution of secondary electrostatic interactions [30], it was decided to keep this term in the scoring function.

Function #7, called SCORE2, reproduces the data with a standard deviation of 7.3 kJ/mol (number of structures n = 82, correlation coefficient r = 0.890, standard deviation s = 7.3, Fisher significance ratio F = 40.3). In this function, $\alpha$ is set to 0.5 and TOL is increased to 0.25 Å. A plot of the calculated $K_i$ values versus the experimentally observed data is shown for function #7 in Figure 1.

Function #7 was further applied to 12 protein–ligand complexes of known three-dimensional structure which were not included in the calibration dataset. This second dataset also contains all structures from our previous dataset [43] that were excluded from the calibration dataset due to their low resolution. It should be noted that this dataset contains six pro-

*Table 2.* Individual contributions, standard deviations s and correlation coefficients r obtained from a fit of free energy functions #1–#8 to experimental binding constants of 82 protein–ligand complexes (all values except r, α and TOL are given in kJ/mol)

| # | 1 | 2 | 3 | 4 | 5 | 6 | 7[a] | 8 |
|---|---|---|---|---|---|---|---|---|
| $\Delta G_0$ | +5.4 | −1.4 | −2.6 | −2.9 | −1.5 | −1.8 | −2.8 | 0.0[b] |
| $\Delta G_{hb}$ | −4.7 | −3.1 | −3.3 | −3.2 | −3.3 | −3.7 | −3.2 | −3.4 |
| $\Delta G_{ionic}$ | −8.3 | −6.6 | −7.0 | −6.1 | −6.1 | −6.2 | −5.7 | −5.9 |
| $\Delta G_{lipo}$ | −0.17 | −0.15 | −0.15 | −0.14 | −0.10 | −0.09 | −0.09 | −0.10 |
| $\Delta G_{rot}$ | +1.4 | +1.0 | +1.0 | +0.9 | +1.1 | +1.1 | +1.0 | +1.1 |
| $\Delta G_{aro}$ | – | – | – | −3.1 | −2.5 | −2.8 | −2.6 | −2.6 |
| $\Delta G_{lipo\ water}$ | – | – | – | – | −1.1 | −1.2 | −1.3 | −1.4 |
| $\Delta G_{esrep}$ | – | – | – | – | – | +0.6 | +0.5 | +0.4 |
| α | 0 | 0 | 0 | 0 | 0 | 0 | 0.5 | 0.5 |
| β | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.2 | 1.2 |
| TOL (Å) | 0.20 | 0.20 | 0.20 | 0.20 | 0.20 | 0.20 | 0.25 | 0.25 |
| s | 9.5 | 8.8 | 8.6 | 8.1 | 7.5 | 7.4 | 7.3 | 7.4 |
| r | 0.835 | 0.841 | 0.837 | 0.859 | 0.882 | 0.887 | 0.890 | 0.890 |

[a] The 95% confidence intervals for function #7 are: $\Delta G_0$: ±158%, $\Delta G_{hb}$: ±29%, $\Delta G_{ionic}$: ±25%, $\Delta G_{lipo}$: ±40%, $\Delta G_{rot}$: ±43%, $\Delta G_{aro}$: ±67%, $\Delta G_{lipo\ water}$: ±48%, $\Delta G_{esrep}$: ±158%.

[b] In function #8 $\Delta G_0$ was set to 0.0 kJ/mol.

teins not present in the calibration dataset. The results are summarized in Table 3. The experimentally determined binding constants [73, 112–121] are predicted by function #7 with a root mean square deviation of 1.54 (log $K_i$) corresponding to an error of 8.8 kJ/mol. A plot of the predicted $K_i$ values versus the experimentally observed data is shown in Figure 2. SCORE2 is now able to predict the binding affinity of the complex antibody–fluoresceine (4FAB). The largest deviation is found for the structure 2PHH. In this complex, the ligand is completely buried in a binding site formed mainly by aromatic side chains.

The constant contribution $\Delta G_0$ in function #7 is negative. In order to assess the influence of $\Delta G_0$ on the quality of the fit, the constant was set to zero in function #8. This leads to a marginal increase of the standard deviation.

## Discussion and conclusions

We have presented an improved set of simple empirical functions that estimate the free energy of binding for a given protein–ligand complex with known 3D structure. The functions were calibrated using a set of 82 protein–ligand complexes with known binding constants. The best representation of the binding data is obtained with function #7 which contains seven
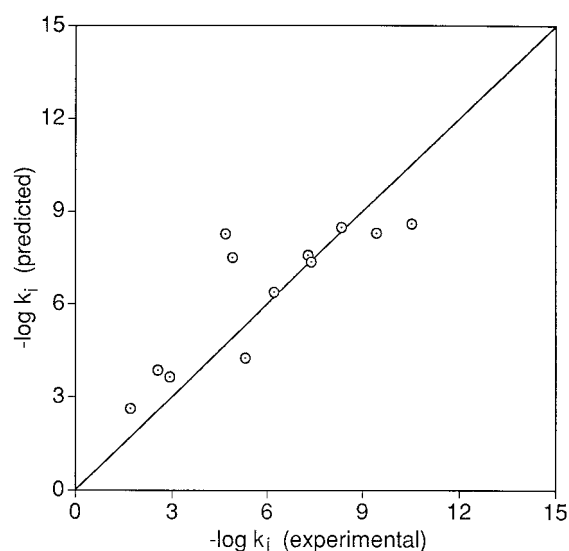


*Figure 2.* Plot of the calculated binding constants $K_i$ from function #7 versus experimentally observed values for 12 protein–ligand complexes not included in the calibration dataset.

adjustable parameters and reproduces the binding energies with a standard deviation of 7.3 kJ/mol. The function was further applied to 12 protein–ligand complexes not included in the calibration dataset. Good agreement is obtained between the predicted and the experimental binding constants. The speed to calculate the new scoring function is comparable to that of

*Table 3.* Protein–ligand complexes used to test the free energy function. The $-\log K_i$ predicted values refer to results obtained with function #7

| Protein–ligand complex | PDB entry | $-\log K_i$ Pred. | $-\log K_i$ Exp. | Ref. |
|---|---|---|---|---|
| Acetylcholinesterase – tacrine | 1ACJ | 7.58 | 7.30 | 111 |
| Carbonic anhydrase – dorzolamide | 1CIL | 8.30 | 9.43 | 112 |
| Thrombin – benzamidine | 1DWB | 3.64 | 2.92 | 113 |
| Thrombin – MQPA | 1DWC | 7.36 | 7.40 | 114 |
| Rhinovirus coat protein – R61837 | 1R09 | 7.50 | 4.90[a] | 115 |
| Purine nucleoside phosphorylase – guanine | 1ULB | 4.25 | 5.30 | 73 |
| PHBH – *p*-hydroxybenzoic acid | 2PHH | 8.27 | 4.68 | 116 |
| Rhinovirus coat protein – Cmpd. IV | 2R04 | 6.38 | 6.22[b] | 117 |
| Antibody – fluorescein | 4FAB | 8.60 | 10.53 | 118 |
| Hemagglutinin – sialic acid | 4HMG | 3.86 | 2.55 | 119 |
| HIV protease – A74704 | 9HVP | 8.48 | 8.35 | 120 |
| FKBP – DMSO | -[c] | 2.62 | 1.70 | 121 |

[a]$IC_{50}$ value.
[b]MIC: minimum inhibitory concentration.
[c]Coordinates were kindly provided by Prof. Walkinshaw.

SCORE1. Therefore, the new function is well suited as scoring function in a 3D database search or de novo ligand design program. The new scoring function SCORE2 is applicable to a broader range of structures with slightly better accuracy than SCORE1.

The values for the ideal neutral hydrogen bond and the ionic interaction are now treated as being dependent on their atomic environment. Within our dataset the contribution from an ideal hydrogen bond (having an ideal geometry) varies between $-1.7$ and $-4.4$ kJ/mol. The contribution from an ionic interaction is between $-3.0$ and $-7.9$ kJ/mol. Both values fall well into the range of the experimental values reported by Fersht [122], Shirley et al. [123] and others [124, 125]. The contribution due to lipophilic contacts $\Delta G_{lipo}$ is predicted to be $-0.10$ kJ/mol $\text{Å}^2$ which is very close to the value estimated by Richards [126] and smaller than the recent estimate of $-0.20$ kJ/mol $\text{Å}^2$ given by Sharp et al. [127]. However, it should be noted that the value for $\Delta G_{lipo}$ is much smaller than in SCORE1 because there are now two additional parameters which also account for apolar interactions. The contribution from rotatable bonds $\Delta G_{rot}$ (+1.1 kJ/mol) is slightly smaller than previous estimates of $+1.6...+3.6$ kJ/mol given by Williams et al. [128]. This low value is probably due to an averaging over protein–ligand complexes with a pre-organized ligand and complexes with a ligand that has to change its conformation upon binding to the protein. The actual value of $-2.6$ kJ/mol obtained

for the interaction between two aromatic rings is also in agreement with estimates by Burley and Petsko ($-2.5...-5.3$ kJ/mol) [48]. The contribution from $\Delta G_{esrep}$ is +0.4 kJ/mol which appears to be physically reasonable. The contribution is essentially a small correction term to hydrogen bonds and ionic interactions.

At first sight it appears that the terms $\Delta G_{lipo}$ and $\Delta G_{lipo\ water}$ describe the same physical effect and are therefore redundant. Nevertheless, a significant improvement of the fit was observed when $\Delta G_{lipo\ water}$ was included in the scoring function. One possible explanation is that a summation of surface contributions does not differentiate between small scattered lipophilic areas and a large continuous lipophilic area. For very small lipophilic areas there exists the possibility that a water molecule is in contact with the lipophilic part of the protein but nevertheless hydrogen-bonded to polar parts of the protein. If a water molecule is in contact with a large continuous lipophilic protein surface, the water molecule cannot form a hydrogen bond with the protein. Therefore, the inclusion of $\Delta G_{lipo\ water}$ offers the possibility to distinguish between these two situations. It should also be noted that $\Delta G_{lipo\ water}$ is a volume dependent term whereas $\Delta G_{lipo}$ is area dependent. The ratio between the volume and the surface of a lipophilic binding site depends on the shape of the binding site. The use of both terms together might offer the possibility to take into account shape effects.

Obviously, the physical interpretation of parameters obtained from a least squares fit to experimental data is not without problems. The limits of a physical interpretation become obvious when looking at $\Delta G_0$. The rationale to include this term was to account for the loss of translational and rotational entropy upon ligand binding. Therefore $\Delta G_0$ should be positive. However, in function #7, a negative value for $\Delta G_0$ was obtained. A physical interpretation of this term is therefore not possible. If $\Delta G_0$ is set to zero, the standard deviation rises by 0.1 kJ/mol with the other parameters essentially left unaffected. The 95% confidence intervals for the terms $\Delta G_0$ and $\Delta G_{esrep}$ are both very large. Therefore, the agreement of the value for $\Delta G_{esrep}$ with results from other work [53] may be fortuitous.

The development of function #7 uses computationally generated positions of discrete water molecules in the binding site. It is well known that most water molecules in binding sites are mobile and are not confined to a single position. In fact, when the calculation of the water positions, as described above, was repeated using different initial water positions, a different water pattern was obtained. The number of unbound water molecules was found to fluctuate by $\pm 1$ for small and tight binding sites such as found for example in trypsin and by $\pm 3$ for large and open binding sites as found in HIV-protease. Nevertheless, in view of the improved accuracy of the scoring function, it was felt that the inclusion of this term is worthwhile. As pointed out above, the contribution from replaced water molecules might be viewed as an additional volume-dependent lipophilic term. Currently, we are investigating alternative approaches for including volume-dependent terms in scoring functions.

There are several limitations to the applicability of SCORE2. The conformational energy of the ligand is not taken into account. The flexibility and conformational change of the protein upon ligand binding are still ignored. It is well known that some proteins undergo drastic conformational changes. It is quite likely that this conformational change yields an increase of the conformational energy, in other words, is detrimental for the binding affinity. On the other hand, it has been shown for several cases that the conformational change of a protein is very similar for the binding of different ligands [129, 130]. Therefore, one can hope that this contribution is similar for different ligands and that the neglect of this contribution does not affect the relative energetic ordering of ligands obtained from the present scoring function.

Another problem is the role of water molecules. It should be noted that SCORE2 can also be used with water molecules that mediate hydrogen bonds between the protein and the ligand. If certain water molecules are thought to be important either structurally or energetically, they can be just added to the protein structure. SCORE2 will then also take into account hydrogen bonds with the water molecules using the same contribution as for interactions with the protein itself.

The dataset used to derive the new scoring function contains protein–ligand complexes where the ligands all show a good steric and electronic complementarity with the protein. The dataset does not contain ligands that do not fit the binding site or show major repulsive electrostatic interactions. The term $\Delta G_{esrep}$ does not differentiate between repulsive secondary interactions as found for example in multiple hydrogen bonds between a ligand and a β-sheet of a protein and strong repulsive electrostatic interactions due to a misplaced polar ligand atom. Therefore, if a ligand forms such a strong repulsive interaction, SCORE2 will overestimate the binding affinity.

Furthermore, it is well known that interactions between a quaternary ammonium group and aromatic rings may contribute significantly to the binding affinity (cation-π interaction) [131]. This effect is still ignored in the current scoring function. However, in view of the very limited amount of quantitative data on this specific interaction, at present it was decided to refrain from incorporating such a term in the scoring function. Interestingly, the binding affinity of the complex acetylcholinesterase–tacrine is predicted very well by SCORE2, although this is thought to be an example where cation-π interaction interactions are important. The problem with very short hydrogen bonds as, for example observed in the complexes thermolysin–ZF$^P$LA and thermolysin–ZG$^P$LL, has been slightly relieved by increasing the tolerances. However, these very short hydrogen bonds are still not recognized as very strong.

Similarly to SCORE1, SCORE2 requires for a successful application a prescreen to detect those ligands that will form some sort of repulsive interaction with the protein, either due to steric problems or due to strong electrostatic repulsions between polar groups. In LUDI [14, 15], this check is carried out in advance for every putative ligand structure and only those structures that pass this test are then subjected to the scoring function. In addition, LUDI can also reject putative ligands that give rise to large voids.

Finally, we note that the present work is based upon experimental data determined in different laboratories. It is clear that the accuracy of any computational approach to predict $\Delta G$ depends on the accuracy of the experimentally determined binding energies. For a number of cases contained in the present dataset several measurements of binding constants were published that show a spread of the experimental values of up to a factor of 5 ($\approx 4$ kJ/mol) [31, 132]. This uncertainty of the experimental binding data poses a limit for the accuracy of any theoretical description of the binding data.

In comparison with other approaches [30–42], the two most notable differences are the use of error tolerances in SCORE2 and the different treatment of solvation and desolvation effects. We refrain from employing an explicit term accounting for solvation. A significant advantage of using such a term would be that calculated solvation energies can be compared with experimental data, at least for small molecules. However, the disadvantage of this approach is that the solvation energies are large and the contribution to $\Delta G$ will be a result from a subtraction of two large numbers. Therefore, the solvation energies have to be calculated with very high accuracy in order to be useful. In the present work, we are exploring an alternative approach to solvation. An implicit description is used. The contribution of a hydrogen bond to the binding affinity is modulated by $f(N_{Neighb})$ and $f_{pcs}$ which account for effects from the molecular neighbourhood of the hydrogen bond. Further, significant desolvation effects are covered by the terms $\Delta G_{lipo}$ and $\Delta G_{lipo\ water}$.

In summary, we have developed a new scoring function with seven adjustable parameters that was calibrated to a dataset of 82 protein–ligand complexes. It is hoped that this new scoring function will increase the predictive power of current computational tools for 3D database searching and de novo design.

## Acknowledgements

## References

1. Lewis, R.A. and Leach, A.R., J. Comput.-Aided Mol. Design, 8 (1994) 467.
2. Nishibata, Y. and Itai, A., J. Med. Chem., 36 (1993) 2921.
3. Bohacek, R.S. and McMartin, C., J. Am. Chem. Soc., 116 (1994) 5560.
4. Gehlhaar, D.K., Moerder, K.E., Zichi, D., Sherman, C.J., Ogden, R.C. and Freer, S.T., J. Med. Chem., 38 (1995) 466.
5. Moon, J.B. and Howe, W.J., Proteins, 11 (1991) 314.
6. Tschinke, V. and Cohen, N.C., J. Med. Chem., 36 (1993) 3863.
7. Rotstein, S.H. and Murcko, M.A., J. Med. Chem., 36 (1993) 1700.
8. Eisen, M.B., Wiley, D.C., Karplus, M. and Hubbard, R.E., Proteins, 19 (1994) 199.
9. Caflish, A., Miranker, A. and Karplus, M., J. Med. Chem., 36 (1993) 2142.
10. Lewis, R.A., Roe, D.C., Huang, C., Ferrin, T.E., Langridge, R. and Kuntz, I.D., J. Mol. Graphics, 10 (1992) 66.
11. Mata, P., Gillet, V.J., Johnson, P., Lampreia, J., Myatt, G.J., Sike, S. and Stebbings, A.L., J. Chem. Inf. Comput. Sci., 35 (1995) 479.
12. Bartlett, P.A., Shea, G.T., Telfer, S.J. and Waterman, S., In Roberts, S.M. (Ed.) Molecular Recognition: Chemical and Biological Problems, Royal Society of London, London, (1989) pp. 182–196.
13. Pearlman, D.A. and Murcko, M.A., J. Comput. Chem., 14 (1993) 1184.
14. Böhm, H.J., J. Comput.-Aided Mol. Design, 6 (1992) 61.
15. Böhm, H.J., J. Comput.-Aided Mol. Design, 6 (1992) 593.
16. Blaney, J.M. and Dixon, J.S., Perspect. Drug Discov. Design, 1 (1993) 301.
17. Kuntz, I.D., Meng, E.C. and Shoichet, B.K., Acc. Chem. Res., 27 (1994) 117.
18. DesJarlais, R.L., Sheridan, R.P., Seibel, G.L., Dixon, J.S., Kuntz, I.D. and Venkataraghavan, R., J. Med. Chem., 31 (1988) 722.
19. Meng, E.C., Shoichet, B.K. and Kuntz, I.D., J. Comput. Chem., 13 (1992) 505.
20. Meng, E.C., Gschwend, D.A., Blaney, J.M. and Kuntz, I.D., Proteins, 17 (1993) 266.
21. Kuntz, I.D., Science, 257 (1992) 1078.
22. Shoichet, B.K., Stroud, R.M., Santi, D.V., Kuntz, I.D. and Perry, K.M., Science, 259 (1993) 1445.
23. Lawrence, M.C. and Davis, P.C., Proteins, 12 (1992) 31.
24. Miller, M.D., Kearsley, S.K., Underwood, D.J. and Sheridan, R.P., J. Comput.-Aided Mol. Design, 8 (1994) 153.
25. Gehlhaar, D.K., Verkhivker, G.M., Rejto, P.A., Sherman, C.J., Fogel, D.B., Fogel, L.J. and Freer, S.T., Chem. Biol., 2 (1995) 317.
26. van Gunsteren, W.F. and Weiner, P.K., Computer Simulations of Biomolecular Systems, ESCOM, Leiden, 1989.
27. a. Kollman, P.A., Chem. Rev., 93 (1993) 2395.
    b. Kollman, P.A., Curr. Opin. Struct. Biol., 4 (1994) 240.
28. Warshel, A., Tao, H., Fothergill, M. and Chu, Z.T., Isr. J. Chem., 34 (1994) 253.
29. Honig, B. and Nicholls, A., Science, 268 (1995) 1144.

30. Grootenhuis, P.D.J. and Van Galen, P.J.M., Acta Crystallogr., D51 (1995) 560.
31. Kurinov, I.V. and Harrison, R.W., Nature Struct. Biol., 1 (1994) 735.
32. Holloway, M.K., Wai, J.M. and Halgren, T.A., J. Med. Chem., 38 (1995) 305.
33. Miranker A. and Karplus, M., Proteins Struct. Funct. Genet., 11 (1991) 29.
34. Vajda, S., Weng, Z., Rosenfeld, R. and DeLisi, C., Biochemistry, 33 (1994) 13977.
35. Krystek, S., Stouch, T. and Novotny, J., J. Mol. Biol., 234 (1993) 661.
36. Horton, N. and Lewis, M., Protein Sci., 1 (1992) 169.
37. Bohacek, R.S. and McMartin, C., J. Med. Chem., 35 (1992) 1671.
38. Jedrzejas, M.J., Singh, S., Brouillette, W.J., Air, G.M. and Luo, M., Proteins, 23 (1995) 264.
39. Nauchitel, V., Villaverde, M.C. and Sussman, F., Protein Sci., 4 (1995) 1356.
40. Peräkylä, M. and Pakkanen, T.A., Proteins, 20 (1994) 367.
41. Head, R.D., Smythe, M.L., Oprea, T.I., Waller, C.L., Green, S.M. and Marshall, G.R., J. Am. Chem. Soc., 118 (1996) 3959.
42. Jain, A.N., J. Comput.-Aided Mol. Design, 10 (1996) 427.
43. Böhm, H.J., J. Comput.-Aided Mol. Design, 8 (1994) 243.
44. Ajay and Murcko M.A., J. Med. Chem., 38 (1995) 4953.
45. Böhm, H.J., J. Comput.-Aided Mol. Design, 8 (1994) 623.
46. Mack, H., Pfeiffer, T., Hornberger, W., Böhm., H.J. and Höffken, H.W., J. Enzyme Inhibition, 9 (1995) 73.
47. Rarey, M., Wefing, S. and Lengauer, T., J. Comput.-Aided Mol. Design, 10 (1996) 41.
48. Burley, S.K. and Petsko, G.A., Science, 229 (1985) 23.
49. Hunter, C.A., Singh, J. and Thornton, J.M., J. Mol. Biol., 218 (1991) 837.
50. Hunter, C.A., Chem. Soc. Rev., (1994) 101.
51. Sander, C., personal communication.
52. GROMOS, user manual.
53. Jorgensen, W.L. and Pranata, J., J. Am. Chem. Soc., 112 (1990) 2008.
54. Connolly, M.L., Science, 221 (1983) 709.
55. Programs INSIGHT and DISCOVER, distributed by MSI, San Diego, CA.
56. Dauber-Osguthorpe, P., Roberts, V.A., Osguthorpe, D.J., Wolff, J., Genest, M. and Hagler, A.T., Proteins, 4 (1988) 31.
57. Kleywegt, G.J. and Jones, T.A., Acta Crystallogr., D50 (1994) 178.
58. Bernstein, F.C., Koetzle, T.F., Williams, G.J.B, Meyer Jr., E.F., Brice, M.D., Rodgers, J.R., Kennard, O., Shimanouchi, T. and Tasumi, T., J. Mol. Biol., 112 (1977) 535.
59. Wilson, D.K. and Quiocho, F.A., Biochem., 32 (1993) 1689.
60. Harvey, S.C., In Goodman, A. (Ed.) The Pharmacological Basics of Therapeutics, MacMillan Press, New York, NY (1987) pp. 980–985.
61. Mangani, S., Carloni, P. and Orioli, P., J. Mol. Biol., 223 (1992) 573.
62. Cappalonga, A.M., Alexander, R.S. and Christianson, D.W., J. Biol. Chem., 267 (1992) 19192.
63. Xiang, S., Short, S.A., Wolfenden, R. and Carter, C.W., Biochemistry, 34 (1995) 4516.
64. Mattos, C., Rasmussen, B., Ding, X., Petsko, G.A. and Ringe, D., Nature Struct. Biol., 1 (1994) 55.
65. Van Duyne, G.D., Standaert, R.F., Karplus, P.A., Schreiber, S.L. and Clardy, J., Science, 252 (1991) 839.
66. Kim, E.E., Baker, C.T., Dwyer, M.D., Murcko, M.A., Rao, B.G., Tung, R.D. and Navia, M.A., J. Am. Chem. Soc., 117 (1995) 1181.
67. Lam, P.Y.S., Jadhav, P.K., Eyermann, C.J., et al., Science, 263 (1994) 380.
68. Morton, A. and Matthews, B.W., Biochemistry, 34 (1995) 8576.
69. White, J.L., et al., J. Mol. Biol., 102 (1976) 759.
70. Bolognesi, M., Cannilo, E., Ascenzi, P., Giacometti, G.M., Merli, A. and Brunori, M., J. Mol. Biol., 158 (1982) 305.
71. Lipscomb, J.D., Biochemistry, 19 (1980) 3590.
72. Turk, D., Stürzebecher, J. and Bode, W., FEBS Lett., 287 (1991) 133.
73. Zollner, H., Handbook of Enzyme Inhibitors, VCH Publishers, Weinheim, 1993.
74. Cowan, S.W., Newcomer, M.E. and Jones, T.A., Proteins, 8 (1990) 44.
75. Wood, J., J. Cardiovasc. Pharm., 14 (1989) 221.
76. Jacobson, B.L., He, J.J., Vermersch, P.S., Lemon, D.D. and Quiocho, F.A., J. Biol. Chem., 266 (1991) 5220.
77. Weber, P.C., Wendoloski, J.J., Pantoliano, M.W. and Salemme, F.R., J. Am. Chem. Soc., 114 (1992) 3197.
78. Matthews, B.W., Acc. Chem. Res., 21 (1988) 333.
79. Fisher, M.T. and Sligar, S.G., J. Am. Chem. Soc., 107 (1985) 5018.
80. Teplyakov, A., Wilson, K.S., Orioli, P. and Mangani, S., Acta Crystallogr., D49 (1993) 534.
81. Cooper, J., Foundling, S., Hemmings, A. and Blundell, T., Eur. J. Biochem., 169 (1987) 215.
82. Miller, D.M., Olson, J.S., Pflugrath, J.W. and Quiocho, F.A., J. Biol. Chem., 258 (1983) 13665.
83. Watson, K.A., Mitchell, E.P., Johnson, L.N., Son, J.C., Bichard, C.J.F., Orchard, M.G., Fleet, G.W.J., Oikonomakos, N.G., Leonidas, D.D., Kontou, M. and Papageorgioui, A., Biochemistry, 33 (1994) 5745.
84. Lowe, J.B., Sacchettini, J.C., Laposata, M., McQuillan, J.J. and Gordon, J.I., J. Biol. Chem., 262 (1987) 5931.
85. Appelt, K., Bacquet, R.J., Bartlett, C.A., et al., J. Med. Chem., 34 (1991) 1925.
86. Bunting, J.W. and Myer, C.D., Can. J. Chem., 53 (1975) 1993.
87. Bolin, J.T., Filman, D.A., Matthews, D.A., Hamlin, R.C. and Kraut, J., J. Biol. Chem., 257 (1982) 13650.
88. Mares-Guia, M. and Shaw, E., J. Biol. Chem., 240 (1965) 1579.
89. Bode, W., J. Mol. Biol., 127 (1979) 357.
90. Wallace, R.A., Kurtz, A.N. and Niemann, C., Biochemistry, 2 (1963) 824.
91. Dani, M., Manca, F. and Rialdi, G., Biochim. Biophys. Acta, 667 (1981) 108.
92. Blaney, J.M., Hansch, C., Silipo, C. and Villon, A., Chem. Rev., 84 (1984) 333.
93. Blundell, T.L., Cooper, J., Foundling, S.I., Jones, D.M., Atrash, B. and Szelke, M., Biochemistry, 26 (1987) 5585.
94. Janes, W. and Schultz, G.E., J. Biol. Chem., 265 (1990) 10443.
95. Miller, M., Schneider, J., Sathyanarayana, B.K., Toth, M.V., Marshall, G.R., Clawson, L., Selk, L., Kent, S.B.H. and Wlodawer, S., Science, 246 (1989) 1149.
96. Bone, R., Vacca, J.P., Anderson, P.S. and Holloway, M.K., J. Am. Chem. Soc., 113 (1991) 9382.
97. Welles, T.N.C. and Fersht, A.R., Biochemistry, 25 (1986) 1881.

98. Verlinde, C.L.M.J., Noble, M.E.M., Kalk, K.H., Groendijk, H., Wierenga, R.K. and Hol, W.G.J., Eur. J. Biochem., 198 (1991) 53.

99. Roderick, S.L., Fournie-Zuliski, M.C., Roques, B.P. and Matthews, B.W., Biochemistry, 28 (1989) 1493.

100. Schloss, J.V., Emptage, M.H. and Cleland, W.W., Biochemistry, 23 (1984) 4572.

101. Kim, H. and Lipscomb, W.N., Biochemistry, 29 (1990) 5546.

102. Lindquist, R.N., Lynn, J.L. and Lienhard, G.E., J. Am. Chem. Soc., 95 (1973) 8762.

103. Kim, H. and Lipscomb, W.N., Biochemistry, 30 (1991) 8171.

104. McPhalen, C.A., Vincent, M.G. and Jansonius, J.N., J. Mol. Biol., 225 (1992) 495.

105. Hol, W. and Verlinde, C., personal communication.

106. Bode, W., Turk, D. and Stürzebecher, J., Eur. J. Biochem., 193 (1990) 175.

107. Stürzebecher, J., Walsmann, P., Voigt, B. and Wagner, G., Thrombosis Res., 36 (1984) 457.

108. Hilpert, K., Ackermann, J., Banner, D.W., Gast, A., Gubernator, K., Hadvary, P., Labler, L., Müller, K., Schmid, G., Tschopp, T.B. and van de Waterbeemd, H., J. Med. Chem., 37 (1994) 3889.

109. Mitchell, E.P., Watson, K.A., Bichard, C., Fleet, G.W.J., Zographos, S.E., Oikonomakos, N.G., Board, M. and Johnson, L.N., In Hunter, W.N., Thornton, J.M. and Bailey, S. (Eds.) Making the Most of your Model, Proceedings of the CCP4 study weekend, Chester, 1995, pp. 111–119.

110. Brandstetter, H., Turk, D., Hoeffken, H.W., Grosse, D., Stürzebecher, J., Martin, P.D., Edwards, B.F.P. and Bode, W., J. Mol. Biol., 226 (1992) 1085.

111. Steinberg, G.M., Mednick, M.L., Maddox, J. and Rice, R., J. Med. Chem., 18 (1975) 1056.

112. Greer, J., Erickson, J.W., Baldwin, J.J. and Varney, M.D., J. Med. Chem., 37 (1994) 1035.

113. Markwardt, F., Walsmann, P. and Landmann, H., Pharmazie, 25 (1970) 551.

114. Kikumoto, R., Tamao, Y., Tezuka, T., Tonomura, S., Hara, H., Ninomiya, K., Hijikata, A. and Okamoto, S., Biochemistry, 23 (1984) 85.

115. Kim, K.H., Willingmann, P., Gong, Z.X., et al., J. Mol. Biol., 230 (1993) 206.

116. Entsch, B., Ballou, D.P. and Massey, V., J. Biol. Chem., 251 (1976) 2550.

117. Badger, J., Minor, I., Kremer, M.J., Oliveira, M.O., Smith, T.J., Griffith, J.P., Guerin, D.M.A., Krishnaswamy, S., Luo, M., Rossmann, M.G., McKinlay, M.A., Diana, G.D., Dutko, F.J., Fancher, M., Rueckert, R.R. and Heinz, B.A., Proc. Natl. Acad. Sci. USA, 85 (1988) 3304.

118. Herron, J.N., He, X., Mason, M.L., Voss, E.W. and Edmundson, A.B., Proteins, 5 (1989) 271.

119. Sauter, N.K., Bednarski, M.D., Wurzburg, B.A., Hanson, J.E., Whitesides, G.M., Skehel, J.J. and Wiley, D.C., Biochemistry, 28 (1989) 8388.

120. Erickson, J., Neidhart, D.J., VanDrie, J., Kempf, D.J., Wang, X.C., Norbeck, D.W., Plattner, J.J., Rittenhouse, J.W., Turon, M., Wideburg, N., Kohlbrenner, W.E., Simmer, R., Helfrich, R., Paul, D.A. and Knigge, M., Science, 249 (1990) 527.

121. Burkhard, P., Kallen, J., Mikol, V. and Walkinshaw, M.D., In Kungl, A.J., Andrew, P.J. and Schreiber, H. (Eds.) Proceedings of the ICSMB95, 1995, pp. 44–60.

122. Fersht, A.R., Shi, J.P., Knill-Jones, J., Lowe, D.M., Wilkinson, A.J., Blow, D.M., Brick, P., Carter, P., Waye, M.M.Y. and Winter, G., Nature, 314 (1985) 235.

123. Shirley, B.A., Stanssens, P., Hahn, U. and Pace, C.N., Biochemistry, 31 (1992) 725.

124. Connelly, P.R., Aldape, R.A., Bruzzese, F.J., Chambers, S.P., Fitzgibbon, M.J., Fleming, M.A., Itoh, S., Livingstone, D.J., Navia, M.A., Thomson, J.A. and Wilson, K.P., Proc. Natl. Acad. Sci. USA, 91 (1994) 1964.

125. Chen, Y.W., Fersht, A.R. and Henrick, K., J. Mol. Biol., 234 (1993) 1158.

126. Richards, F.M., Annu. Rev. Biophys. Bioeng., 6 (1977) 151.

127. Sharp, K.A., Nicholls, A., Friedman, R. and Honig, B., Biochemistry, 30 (1991) 9686.

128. Searle, M.S., Williams, D.H. and Gerhard, U., J. Am. Chem. Soc., 114 (1992) 10697.

129. Sali, A., Veerapandiam, B., Cooper, J.B., Moss, J.B., Hofmann, T. and Blundell, T.L., Proteins, 12 (1992) 158.

130. Wierenga, R.K., Noble, M.E.M. and Davenport, R.C., J. Mol. Biol., 224 (1992) 1115.

131. Dougherty, D.A. and Stauffer, D.A., Science, 250 (1990) 1558.

132. Baker, B.R. and Erickson, E.H., J. Med. Chem., 10 (1967) 1123.