# *Q-fit*: A probabilistic method for docking molecular fragments by sampling low energy conformational space

Richard M. Jackson*
*Department of Biochemistry and Molecular Biology, University College, Gower Street, London WC1E 6BT, UK;
*Current address: School of Biochemistry and Molecular Biology, University of Leeds, Leeds LS2 9JT, UK
(E-mail: jackson@bmb.leeds.ac.uk)*

## Summary

A new method is presented that docks molecular fragments to a rigid protein receptor. It uses a probabilistic procedure based on statistical thermodynamic principles to place ligand atom triplets at the lowest energy sites. The probabilistic method ranks receptor binding modes so that the lowest energy ones are sampled first. This allows constraints to be introduced to limit the depth of the search leading to a computationally efficient method of sampling low energy conformational space. This is combined with energy minimization of the initial fragment placement to arrive at a low energy conformation for the molecular fragment. Two different search methods are tested involving (i) geometric hashing and (ii) pose clustering methods. Ten molecular fragments were docked that have commonly been used to test docking methods. The success rate was 8/10 and 10/10 for generating a close solution ranked first using the two different sampling procedures. In general, all five of the top ranked solutions reproduce the observed binding mode, which increases confidence in the predictions. A set of ten molecular fragments that have previously been identified as problematic were docked. Success was achieved in 3/10 and 4/10 using the two different methods. Again there is a high level of agreement between the two methods and again in the successful cases the top ranked solutions are correct whilst in the case of the failures none are. The geometric hashing and pose clustering methods are fast averaging $\sim 13$ and $\sim 11$ s per placement respectively using conservative parameters. The results are very encouraging and will facilitate the process of finding novel small molecule lead compounds by virtual screening of chemical databases.

## Introduction

Computer-aided strategies have become an increasingly important component of structure-based drug design. The aim is to identify small molecules that bind and inhibit or otherwise alter activity of the biological target. Current computational methods involve database screening of small molecule libraries by molecular docking [1–3] or the *de novo* design of ligands by assembly of molecules from smaller fragments that interact favourably with the target [4–7]. The underlying complexity of calculating non-covalent interactions in solution and the enormous chemical diversity of potential small molecule drugs means that the ability to develop more accurate and time efficient methods is an active area of research.

The prediction of the free energy of binding in biological systems is still poorly developed. However, the use of computational methods such as docking and *de novo* design by fragment building are of significant value as the number of possible compounds in molecular space is too large to screen experimentally. Design strategies require a manageable subset even with experimental high-throughput screening methods [8]. In one study virtual screening provided lead compounds where conventional random screening had failed [37]. The requirement for time efficient search and scoring procedures has therefore become of in-

creasing focus due to the enormity of small molecule databases. There are currently several docking methods that focus on the accurate prediction of the ligand binding mode (where the time required is a secondary concern). These include methods that employ stochastic sampling techniques including simulated annealing [9–10], genetic algorithms [3, 11], evolutionary programming [12] or a Tabu search methodology [13]. Alternatively, deterministic approaches limit the sampling space to arrive at solutions in a faster time scale, these include DOCK [1, 36], FlexX [2] and CLIX [23] as well as several other algorithms (for a review see [14]).

Comparing the accuracy of the different methods is an ongoing issue [15], however, speed is an important consideration for application to drug design. The virtual screening of a large compound library of 100 000 compounds remains impractical for all but the fastest methods using multiprocessor computing. Spending only one minute per compound would take about ten weeks to search such a library on a single processor machine. Alternatively *de novo* design assembles ligands from small fragments that are identified to bind favourably. These methods originate from those used to model the interaction of functional groups with a protein such as Goodford's GRID program [16]. Examples of these methods include the Multiple copy simultaneous search method, MCSS [4] and HOOK [7] which connects MCSS generated fragments, as well as LUDI [5], GROW [17], SPROUT [6] and others (for a review of methods see [18] ). The degree of chemical diversity that can be explored is considerable, the limitations are that they produce ligands that are not always accessible synthetically and therefore not necessarily commercially viable.

This article addresses the problem of accurate prediction of the binding mode of a ligand to a receptor a problem commonly referred to as 'the docking problem'. The problem of correctly predicting binding affinity or correctly identifying active molecules from a large database has not been addressed here. Emphasis has been placed on docking rigid fragments to a receptor. The inclusion of ligand and protein flexibility [19] will be addressed in a subsequent study. The method involves using a new probabilistic procedure for determining positions and orientations of small sized molecular fragments in the binding site of a rigid protein. The aim is to facilitate the process of finding novel small molecule lead compounds that bind in a specific way to the receptor. Potentially these small molecule fragments could be connected using com-

binatorial principles to create *de novo* compounds or used as the base fragment for searching existing chemical databases from which matches could be docked using flexible ligand and protein models. The program uses a probabilistic method to place ligand atom triplets at the lowest energy receptor sites. The use of geometric hashing [21, 25] and pose clustering [2] have been used in other molecular docking programs and probabilistic ranking of receptor binding modes has also been considered independently [23]. This algorithm combines probabilistic ranking with these powerful search methods. Previous implementations [2, 25] do not use probabilistic ranking, but simply search all defined ligand triplets for compatiblity with all receptor triplets of similar probe types. Lawrence and Davis [23] define binding modes by receptor probe pairs, followed by a rotational search. Here the binding mode is uniquely defined in 3D in terms of a triplet probe interaction. The ranking of triplet energies allows an ordered search of binding modes testing the most favourable triplet interactions first. Therefore, constraints can be placed on the depth of the search whilst guaranteeing that the lowest energy triplet interactions are sampled first. Once the fragment is placed in the binding site grid-based energy minimization is performed using the simplex method.

**Theory and methods**

*A probabilistic model for binding*

From a statistical thermodynamic viewpoint the energies of atoms and therefore molecules are confined to discrete values. Therefore, the allowed energy levels can in principle be used to describe the behaviour of bulk matter. From a knowledge of these energy levels at a single instant we can calculate the average of the property from all members of an ensemble and this is equivalent (the *ergodic* hypothesis) to solving the more difficult problem of finding the time average for a single system. One of the outcomes of this statistical description is the Boltzmann distribution describing the most probable configuration of an ensemble by;

$$P(E_i) = n_i/N = \exp(-E_i/kT)/\Sigma \exp(-E_i/kT), \quad (1)$$

Where the sum over the states of the members of the ensemble, $\Sigma \exp(-E_i/kT)$, is the *partition function, Q*. In terms of a system involving the free energy for the interaction of a single functional group or probe

(confined to a grid) with a protein active site, the solution to finding the most probable configuration is trivial. It is given by the minimum energy interaction of the probe on the grid. For two connected (bonded) point probes the solution becomes more complex, and we must consider the probe types as well as the distance between the two probes. However, given two different probes with associated grid energies, we can sort the grid points in both grid types in order of decreasing favourable interaction energy. Given a finite number of points the optimal pair and an ordered list of all sub-optimal pairs of points can be constructed by further sorting pairs in terms of overall interaction energy. In order to find the lowest energy interaction we must sort down through this list of pairs until we find one that satisfies the conditions (i) that both types of probe are present and (ii) the known distance constraint between the bonded probes is satisfied. Given that we can describe the system by the grid/probe approximation this match will be the interaction energy *global* minimum for the system. This treatment can be extended to three-, four-, through to *n*-point probe interactions although at ever increasing computational expense.

In the current application of the program we consider only three-point probe interactions (see below), as a molecular fragment transformation is uniquely defined by mapping three probe points onto three receptor interaction points [2] and computational overheads are kept to a minimum. As ligand fragments are generally defined by more than three probe centres we can no longer define the *global* minimum in terms of a three point interaction. As with molecular energies there is also predicted to be a Boltzmann-like distribution (Equation 1) of individual atom (probe) energies. Therefore the three most favourable probe interactions will statistically constitute a considerable proportion of the overall molecular interaction energy. In such cases a very limited number of mappings may therefore be considered whilst still expecting to see the interaction energy *global* minimum represented.

### Probabilistic docking and the Q-fit program

Q-fit is a suit of programs that (1) Adds hydrogen atoms and lone pair electrons to existing protein PDB [20] coordinates. (2) Calculates the interaction energy of a user defined list of atom probes (representing the ligand or series of ligands) with the protein receptor using the GRID force field [16] and stores them in 3D grid format. (3) Docks the coordinates of the ligand to the receptor binding site for which interaction energies have already been calculated in (2). Interaction grids are only calculated once for a given target receptor. The steps are described below in more detail. Firstly the docking algorithm is described; then the empirical energy function and calculation of the non-bonded interaction energy; then the placement of the hydrogen atoms and lone pairs. Finally, the initial startup conditions of the receptor-ligand systems and program parameters are detailed.

### Defining the most energetically favourable binding site locations

Figure 1 is a schematic description of the steps involved in the docking method described here. Given the pre-calculated 3D grid maps of atom preferences calculated in (2) the program reads those appropriate for the atom probe types in the given ligand. For each map the grid point interaction energies are then sorted in order of decreasing interaction energy with the receptor so that the most favourable (i.e. negative) interaction energy is at the top of the list and the least favourable at the bottom (Figure 1A). For each probe type the top, $N$, locations are found such that they are not within a defined cut-off distance, $R_{cut}$, of one another. Both $N$ and $R_{cut}$ are user defined parameters (see below) which determine the potential depth of the search and geometric proximity of the minima respectively. This list of $N$ sites for each probe type are combined and again sorted in order of decreasing interaction energy irrespective of the probe type. This set of, $M$, interaction sites ($N^*$ number of atom types in the ligand) corresponds to a ranked list of the most energetically favourable binding sites.

### Matching ligand atom triplets to those of the receptor target

The aim of the algorithm is to match a triplet corresponding to the ligand atoms (Figure 1C) to those of the most favourable binding location in the receptor (Figure 1D). Two alternate algorithms from computer vision object recognition have been tested for suitability to perform this task.

### Search involving a geometric hashing method

This procedure uses geometric hashing [21] and involves histograming (or binning) the ligand atom-atom distances for each possible triplet combination. The differences in distances are separated into bins of size,
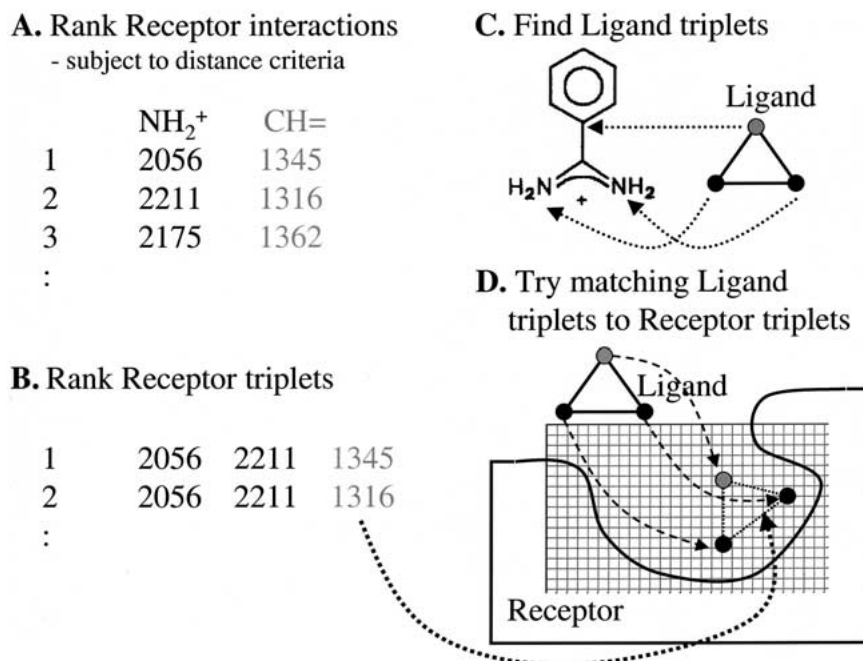
**A. Rank Receptor interactions**
- subject to distance criteria

|   | $NH_2^+$ | $CH=$ |
|---|---|---|
| 1 | 2056 | 1345 |
| 2 | 2211 | 1316 |
| 3 | 2175 | 1362 |
| : | | |

**B.** Rank Receptor triplets

| 1 | 2056 | 2211 | 1345 |
|---|---|---|---|
| 2 | 2056 | 2211 | 1316 |
| : | | | |

**C. Find Ligand triplets**

**D.** Try matching Ligand triplets to Receptor triplets

*Figure 1.* Schematic representation of the steps (A–D) involved in the docking algorithm (see text for details). The numbers in A and B refer to indices of the grid points.

*binsize*. These three integer distances are used as an address to a hash table which includes the identity of the corresponding atom triplet. The ligand hash table is small and compiled at run time. Similarly for the list of $M$ receptor interaction sites the grid-grid distances for each triplet combination are binned such that they could potentially correspond to a triplet in the ligand in terms of atom types (e.g. a triplet with two grid points with the same atom type would be rejected if the ligand only contained one atom of this type). All possible triplets in the receptor are then sorted by energy so that the triplet with the highest combined interaction energy is at the top and lowest at the bottom (Figure 1B). The matching procedure or recognition phase then proceeds taking the top scoring receptor triplet and uses the distances to compute the ligand hash table address(es) that could contain a match. For each ligand triplet stored at that address the atom types are compared to the receptor triplet and matches are stored for the transformation stage (see below).

*Search involving a pose clustering method*

This procedure is based on a pose clustering [2, 22]. The theory implies that we can achieve accuracy equivalent to examining all group matches (3 points matches) by examining sub-problems in which only those group matches that share some basis of two point matches are considered. Normally recognition in computer vision involves choosing a random triplet of image (receptor) points and then all matches with the object are examined and if one leads to recognition of the object we may stop. Otherwise we continue choosing image basis points at random. In molecular docking we do not have such a termination criterion, i.e. we do not know when recognition has been achieved. However, the most energetically favourable interactions are most likely to be found by examining the pairs that give the greatest interaction energy, *i.e.* a probabilistic sampling. The $M$, receptor interaction sites, $\mathbf{m}_i$, corresponding to the ranked list of the most energetically favourable binding sites are examined in a systematic *breadth-first* search (for all $\mathbf{m}_i | i = 1, \ldots, n - 2$ and for all $\mathbf{m}_j | j = i + 1, \ldots, n - 1$) to find a pair ($\mathbf{m}_i$, $\mathbf{m}_j$) that is coincident to a ligand pair ($\mathbf{l}_i$, $\mathbf{l}_j$) such that the distances ($|\mathbf{m}_i - \mathbf{m}_j| - |\mathbf{l}_i - \mathbf{l}_j| < \delta$). Then find a further point match (for all $\mathbf{m}_k | k = 1, \ldots, n$ and $k \neq i$, $k \neq j$) that has the most favourable triplet interaction energy and is also coincident to the ligand triplet ($\mathbf{l}_i, \mathbf{l}_j, \mathbf{l}_k$) such that the distances ($|\mathbf{m}_i - \mathbf{m}_k| - |\mathbf{l}_i - \mathbf{l}_k| < \delta$ and ($|\mathbf{m}_j - \mathbf{m}_k| - |\mathbf{l}_j - \mathbf{l}_k| < \delta$). Such triangle matches are $\delta$-compatible and are stored for the transforma-

tion stage. This method is similar to the method of Lawrence and Davis [23], except a triplet is used as the group match (as opposed to a pair, followed by rotational sampling). We have many more potential receptor (image) matching sites ($N = 200$ per ligand probe type) as opposed to anything between one and thirteen per ligand probe type as in their docking study on hemagglutinin.

In both the geometric hashing and pose clustering algorithms we find optimal triplet interaction sites of the ligand in the receptor using different search methodologies. In both cases a decision must be made about how far down the list we need to search in order to be highly likely to find a solution that approaches the global minimum energy conformation. This value, $L$, is a trade-off between finding the global minimum energy conformation and the run time for a particular ligand-receptor pair. Furthermore, both the lists of candidate solutions can contain identical solutions. Transformations that have identical transformation matrices are checked for and are eliminated from further study. Additionally, solutions that are very similar to ones already generated are also eliminated. In the current implementation if a newly generated solution is within 0.75 Å RMSD of any previous solution it is eliminated from further study. The remaining candidate solutions are transformed using a least squares fitting routine [24] where the ligand atoms are transformed onto the receptor interaction sites. Following placement an atom overlap test is performed between ligand and receptor atoms and only candidate solutions with less than a given percentage, $P_{bumps}$, of atom-probe overlaps or 'bumps' (calculated as a percentage of the number of atoms in the ligand) are retained for further study.

### Energy minimization of the ligand fragment in the receptor

In many docking algorithms little or no energy refinement is carried out following initial placement of the ligand in the receptor. This may be adequate for many empirical energy functions [2, 25] but is unsuitable for use with the molecular mechanics energy function used here. Following steric acceptance of the ligand fragment rigid-body minimization was performed. The downhill Simplex algorithm of Nelder and Mead [26] was implemented according to Gschwend and Kuntz [27]. The method requires only function evaluations and no derivatives. It was found to be considerably faster than a gradient descents minimization

routine implemented previously [19], with results of comparable quality. In the current implementation minimization is deemed complete when a restart simplex reduces the energy by less than 0.01 kcal/mol or a maximum number of 500 iterations exceeded.

### Evaluation of the grid-based interaction energy

Q-fit calculates the ligand-receptor non-bonded interaction energy using force field parameters from GRID (version 15). The non-bonded interaction energy of a probe type with the receptor is calculated at each xyz position on a user defined rectangular grid. The empirical energy function consists of van der Waals (modeled by the 12–6 Lennard–Jones equation), electrostatic and hydrogen bond functions. The van der Waals and electrostatic functions are modelled in an identical way to GRID (for details see [16, 28]). The hydrogen bond function uses a direction-dependent 8-6 function consistent with the hydrogen bond parameters of GRID (version 15) and are calculated according to:

$$E_{hb} = [C_{ij}/d^8 - D_{ij}/d^6]\cos^m \theta, \qquad (2)$$

where

$$C_{ii} = -3E_{min}(2R_{min})^8,$$

$$D_{ii} = -4E_{min}(2R_{min})^6$$

and $E_{min}$ is the minimum in the potential energy well when two identical atoms of type, i, are interacting and $R_{min}$ is half the distance between the atoms at this point. Combining rules for non-identical atoms are geometric for well depth and arithmetic for contact distance. If the receptor donates *or* accepts a hydrogen bond the energy is angle dependent. In the case of the receptor donating a hydrogen bond to a probe, $\theta$, is the angle DHP where, D, is the protein atom donor, H, the hydrogen atom on the donor and P is the probe. In the case of the receptor accepting a hydrogen bond, $\theta$, is the angle ALP where, A, is the protein acceptor, L, is the lone pair of electrons on the acceptor and P is the probe. The hydrogen atoms and lone pair electrons on the protein receptor are computed from the heavy atom coordinates of the protein (see below). It is assumed that the probe hydrogens or lone pairs can orient themselves to optimise their interaction with the protein receptor. The term, $m$, is 4 and $E_{hb}$ is set to zero when $\theta \leq 90°$. The maximum number of hydrogen bonds that can be accepted and donated is specified in the GRID parameters file for the protein

and are also specified for the probe. Of the many possible hydrogen bonds that the probe can make only the energetically most favourable ones are selected subject to this constraint. Overall the main difference from Goodford [16] is inclusion of the angle dependence of protein hydrogen bond acceptors, although this feature has been included in the latest versions of GRID (P. Goodford, pers. commun.).

*Addition of hydrogen atoms and lone pair electrons*

Hydrogen atoms and lone pair electrons were added using existing bonds and geometries with hydrogen atoms being added according to their equilibrium bond lengths, bond angles and dihedral angles [19]. This procedure was adapted to additionally include lone pairs. These were added to acceptors at a distance of 1.0 Å purely for the purpose of calculating the angular dependence of acceptor hydrogen bonds to a probe. They have no van der Waals or electrostatic properties themselves they serve only as dummy atoms to define the hydrogen bond geometry. The lone pair geometries of the $sp^2$ hybridized atoms: carbonyl oxygen, phenol oxygen, imidazole nitrogen, pyridine nitrogen, are trigonal planar. The lone pair geometries of $sp^3$ hybridized atoms: water oxygen, aliphatic alcohol oxygen, ether oxygen, primary and secondary amine nitrogen, are tetrahedral. Note, that the sulphydryl sulphur has no hydrogen bonding capacity in the GRID force field. Multiple hydrogen and lone pair locations are generated in the case of Ser, Thr, Cys, Tyr and His and can interact with the probe simultaneously. For Ser, Thr and Cys(h) the *gauche-, trans* and *gauche+* conformations of the $\chi_2$ dihedral defined as (Cα-Cβ-O(S)γ-H(Lp)γ) are generated. For Tyr the two *cis* and *trans* conformations of the $\chi_6$ dihedral defined as ($C^{\epsilon 1}$-$C^{\xi}$-$O^{\eta}$-H(Lp)$^{\eta}$) are generated. Also the for histidine both the $N^{\epsilon 2}$ and $N^{\delta 1}$ protonated tautomers are generated. Thus, for any interaction of the probe with these potential hydrogen bond donors/acceptors all potential interactions are calculated. Therefore the groups can satisfy their intermolecular hydrogen bonding commitments of any of their conformers (or tautomers in the case of histidine). However, only the most favourable interaction can (potentially) be used. This depends on the constraints for the maximum number of hydrogen bonds that can be accepted/donated by the probe at that particular location.

*Input data and program parameters*

The protein-ligand coordinate files were taken from the PDB database [20] for the appropriate complex. In the case of the receptor protein coordinates hydrogen atoms and lone pairs were added to all standard amino acids and their N and C termini (note: N-terminal proline is a special case) as described above. This is followed by force field parameter assignment of protein atoms and the heavy atoms of hetero atoms and non-standard amino acids for which parameters are assigned in the GRID force field file. Each atom of each residue in the force filed parameters file is checked for missing atoms (alternate locations are ignored) and overall residue charge assignments and missing atom records are recorded. In addition the all atom parameter assignments including added hydrogen atoms and lone pairs are recorded with assigned van der Waals, electrostatic and hydrogen bonding parameters in an extended PDB format.

Probe atoms are defined by the same force field description as the protein. The probe types are included in an input file for the generation of 3D grid maps. The map is produced for each probe type and may constitute only those atom probe types present in a particular ligand or may include all possible probe types where the docking of multiple ligands is to take place. The ligand (or ligand fragment) was extracted from the PDB file format. Ligand atoms in the Cartesian coordinate file are assigned probe atom types for the docking phase. For the ligand representation, only the heavy atoms are defined and the receptor ligand interaction scored according to the grid based energies.

Q-fit requires the definition of the receptor site as a rectangular box. With a user defined box center at $x$, $y$, $z$ and $xyz$ box dimensions d$x$, d$y$, d$z$. In all runs described here the box center was taken as the approximate ligand center and the box dimensions uniformly described by a 20 Å cube. The grid spacing for sampling the receptor-probe interaction can be changed from the default (0.5 Å/grid). Also under user control are: (1) the maximum positive (unfavourable) grid energy (default = 5.0 kcal/mol); (2) the van der Waals non-bonded cut-off distance (default = 8.0 Å); (3) the hydrogen bond cut-off distance (default = 5.0 Å); (4) the electrostatic cut-off distance (default = 12.0 Å); (5) the dielectric constant of the protein interior (default = 4.0); (6) the dielectric constant of water (default = 80.0); (7) the distance defining unacceptably close van der Waals contacts of the probe with polar atoms (default = 2.3 Å) and non-polar atoms

(default = 2.8 Å) for definition of the overlaps or 'bumps' grid (hydrogen atoms and lone pairs are not included in this analysis).

Several parameters must be assigned for the docking stage (see above). Alternative combinations for the values were tested. Here, we report only those used in the results section. They represent an optimized or recommended parameter set. These are: the top, $N$, locations for the each probe type (default = 200); the top, $M$, locations for all probes (default = No. of probe types * $N$); the top, $L$, accepted triplet matches that are considered for scoring (default = 3000); the cut-off distance, $R_{cut}$, between minima (default = 1.0 Å) used in defining interaction sites in both algorithms. The value of binsize for the ligand atom-atom, and interaction site grid-grid distances (default = 3.0 Å), used only in the geometric hashing algorithm. The value of δ the maximum difference in distance between matched ligand atom-atom, and interaction site grid-grid distances (default = 3.0 Å), used only in the pose clustering algorithm. The percentage value, $P_{bumps}$, of allowed atom-probe overlaps or bump, calculated as a percentage of the number of atoms in the ligand (default = 33.0%).

*Program user interface*

Execution of the programs is command-line driven. The user has control over most of the analysis parameters with optional command line arguments. There are three distinct program steps: (1) Generation of hydrogen atoms and lone pair electrons for an existing protein PDB coordinate file. (2) Calculation of the receptor-probe interaction energy 3D grids (for each probe type in a user defined list). (3) Docking using either the pose clustering algorithm or geometric hashing algorithm. It should be noted that stages (1) and (2) need only be performed once for the docking of multiple ligands to a given receptor.

**Results**

The twenty protein-ligand complexes in the dataset were chosen to reflect a diverse set of fragments (see figure 2) that bind non-covalently to their receptors but have few or no rotatable bonds. The ligands are either rigid or are one of a limited number of accessible conformers and fulfil the criteria that justifies rigid-body fragment docking. Additionally, the fragments were chosen to correspond with those of other docking studies to allow comparison. The study by Jones [3] tested

a dataset of 100 complexes from the Protein Data Bank. Ten protein-ligand complexes are chosen from their *Good* or *Close* predictions (hereafter referred to as dataset 1). They are complexes that have also been commonly used in bench-marking other docking methods [2, 13, 29]. Ten protein-ligand complexes were also choosen from their *Errors* or *Wrong* prediction categories to test the algorithm on the most problematic complexes identified in their study (hereafter referred to as dataset 2). Before performing docking all water molecules are removed whilst other bound ligands or cofactors are retained (if present). Note that all RMSD values quoted are for heavy atoms only and are calculated from the experimental crystallographic structure and not a pre-minimized protein-ligand complex. Therefore, the energy minimized experimental structures have a non-zero RMSD value. Any solution with an RMSD of less than 2 Å was deemed to be an acceptably good prediction of the binding mode. Here only one simulation is performed per receptor-ligand pair unlike stochastic sampling methods where many simulations are performed to ensure most of the high-affinity binding modes are explored. Timings are therefore reported for each docking simulation (using the parameters reported in theory and methods section) on a Linux workstation (600 MHz processor). The time used to generate hydrogens/lone pairs on the protein receptor and to perform the grid based energy calculations for different receptor probes were less that 10 s in total for any of the systems reported here. This only needs to be calculated once for a receptor when docking multiple ligand fragments.

*Prediction of the binding mode*

The results of docking the two different ligand data sets to their protein receptors are given in Table 1 for the geometric hashing and Table 2 for the pose clustering algorithms. Solutions are ranked according to predicted receptor-ligand interaction energy. The summary of the results include the RMSD and calculated interaction energies of (1) the top ranked solution (2) the best prediction with an RMSD of less than 2 Å, including its rank and (3) the energy minimized crystallographic solution including its rank. Furthermore, the number of predictions generated in the top five ranked solutions with an RMSD of less than 2 Å is shown, lastly, the elapsed time is given.

The level of success (determined by a correct prediction of the binding mode ranked number one) is high for dataset 1 for both the geometric hashing and
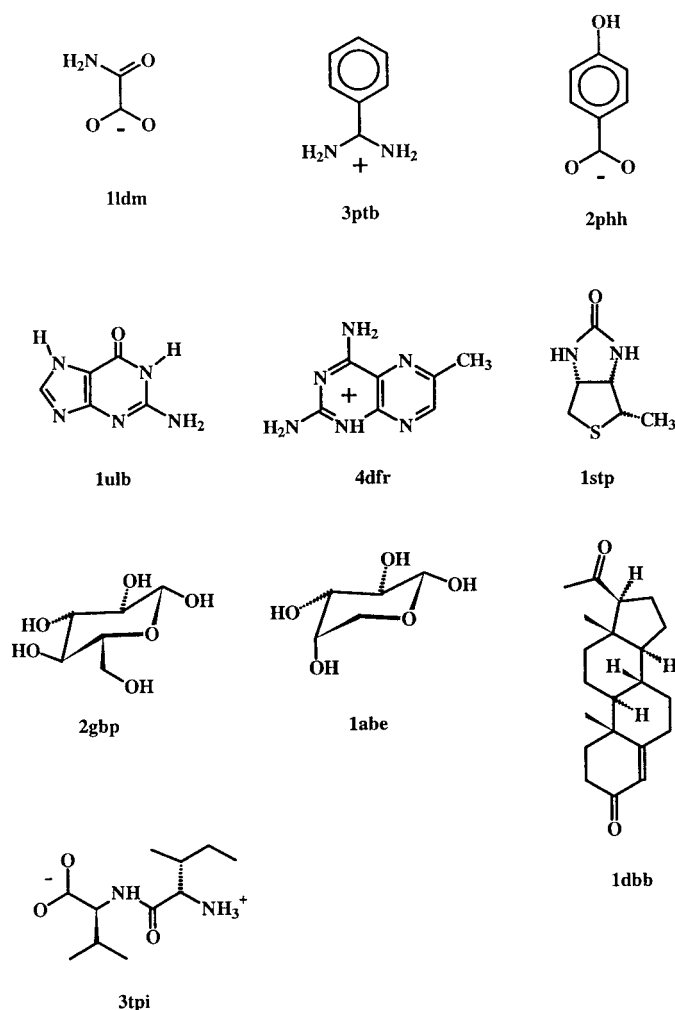
*Figure 2.* Chemical structures of the ligand fragments and the PDB code for the receptor coordinates. 2A dataset 1.

pose clustering algorithms. Eighteen of the twenty independent docking runs correctly predict the binding mode of the ligand ranked first. Similarly the number of predictions with an RMSD of less than 2 Å in the top five solutions is high in both cases. As expected in the case of dataset 2 both algorithms do worse overall. The degree of correlation between the successes and failures of the two algorithms is high. Both algorithms tend to do either well or poorly depending on the particular complex concerned. It would appear that use of either search strategy is valid.

*Analysis of predictions in dataset 1*

The predicted top ranked ligand conformations in dataset 1 are in good agreement with experiment and whilst minor differences are observed they generally have a similar RMSD to the minimized experimental structure. Figure 3 shows a summary of the results for the docking of purine-nucleoside phosphorylase and guanine using the geometric hashing algorithm with a plot of interaction energy versus RMSD for generated placements following energy minimization. The minimized X-ray conformation is ranked number five. All the generated placements ranked in the top ten have an RMSD of less than 1 Å. The graphic of the active site (with the protein represented by a Cα trace) shows the good agreement of the ligand following energy minimization of the experimental X-ray structure (red) with the top four scoring placements. This result is typical of ligand-receptor complexes analysed in dataset 1.

A comparison of Tables 1 and 2 shows that the minimized experimental structures are identical (same
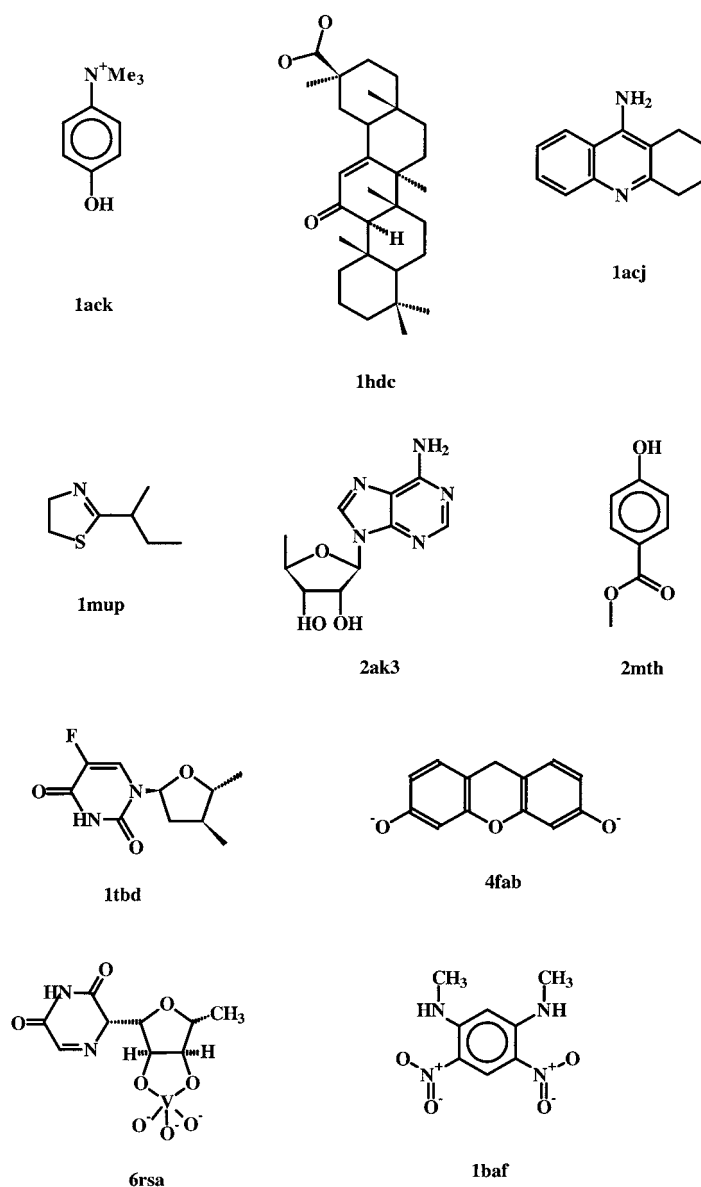
*Figure 2B.* dataset 2.

RMSD and energy) using either of the two algorithms. This is due to the fact that the initiation seed number in the random number generator of the simplex minimization algorithm is identical. Therefore, identical placements of the ligand will lead to the same energy minimum. With the exception of 1dbb (using the geometric hashing algorithm) all the complexes in dataset 1 have more than one solution ranked in the top five with a RMSD (and non-bonded contacts) that is highly similar to the experimental crystallographic conformation. However, some slight variations in binding conformation are worthy of note in the cases of 3ptb.

In the case of the ligand benzamidine binding to trypsin the experimental (3ptb [30]) and predicted complexes (RMSD = 1.474 Å and 0.896 Å in the geometric hashing and pose clustering algorithms respectively) are in a slightly different conformation. The two binding modes are very similar in most respects and have an identical hydrogen bonding pattern with Asp 189, Ser 190 and Gly 219, however, the hydrogen bond lengths are different. The distances be-

# Purine-nucleoside phosphorylase + guanine



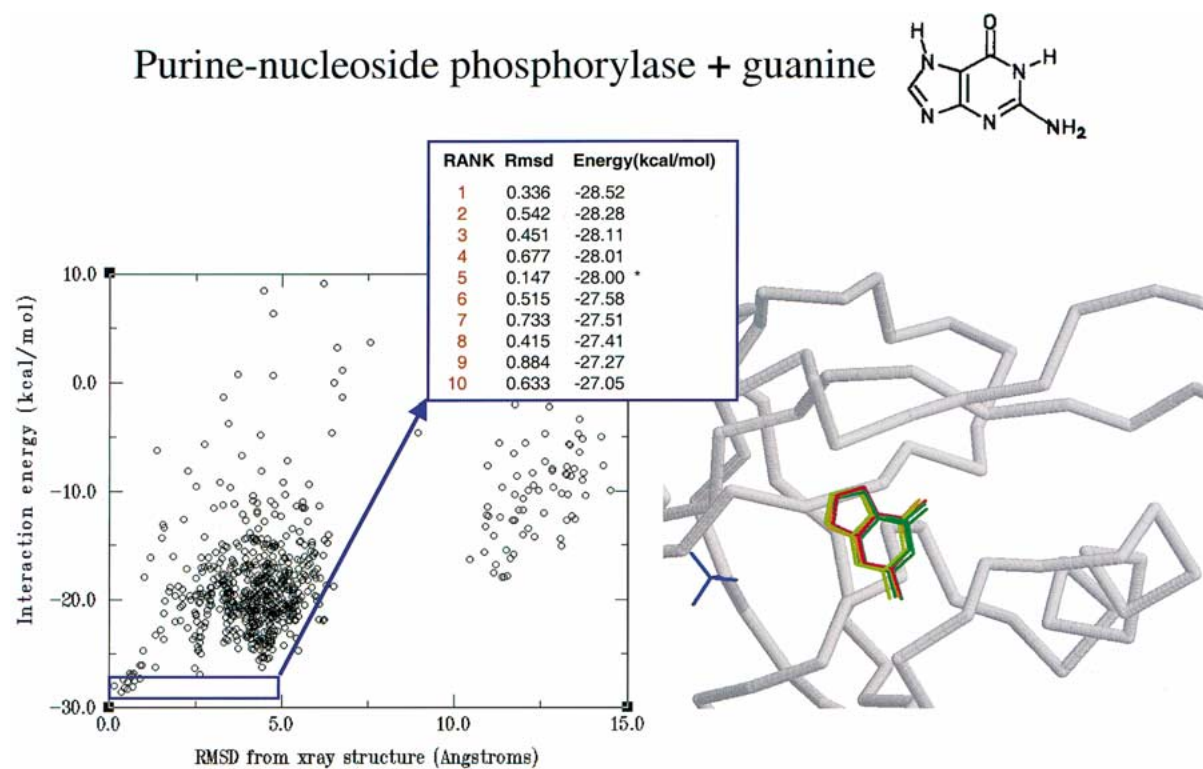| RANK | Rmsd | Energy(kcal/mol) |
|------|------|------------------|
| 1 | 0.336 | -28.52 |
| 2 | 0.542 | -28.28 |
| 3 | 0.451 | -28.11 |
| 4 | 0.677 | -28.01 |
| 5 | 0.147 | -28.00 * |
| 6 | 0.515 | -27.58 |
| 7 | 0.733 | -27.51 |
| 8 | 0.415 | -27.41 |
| 9 | 0.884 | -27.27 |
| 10 | 0.633 | -27.05 |

*Figure 3.* Docking of purine-nucleoside phosphorylase and guanine using the geometric hashing algorithm. See text for details.
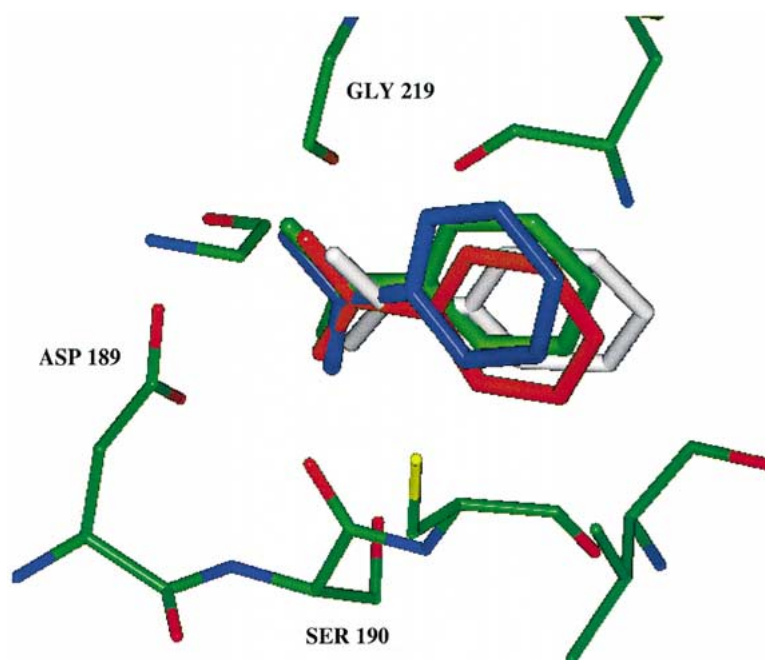


*Figure 4.* Trypsin-benzamidine binding site: experimental X-ray structure (white) and the top three predicted conformations using the geometric hashing algorithm with the top ranked solution in blue.

*Table 1.* Docking of Datasets 1 and 2 using the geometric hashing algorithm

| PDB complex | No. of solutions | Lowest energy[a] | | Best prediction | | | Minimized experimental | | | No. < 2 Å in Top 5 | Elapsed time[c] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | RMSD | $\Delta\Delta E_{LE\rightarrow BP}$[b] | RMSD | $\Delta E$ | Rank | RMSD | $\Delta E$ | Rank | | |
| 1ldm | 255 | 2.415 | 0.6 | 0.788 | −34.4 | 2 | 0.424 | −34.1 | 3 | 3 | 0:06 |
| 3ptb | 167 | 1.474 | 0.0 | 1.474 | −37.7 | 1 | 0.347 | −32.5 | 76 | 5 | 0:06 |
| 2phh | 524 | 0.634 | 0.0 | 0.634 | −32.0 | 1 | 0.375 | −30.0 | 10 | 4 | 0:10 |
| 1ulb | 694 | 0.336 | 0.0 | 0.336 | −28.5 | 1 | 0.147 | −28.0 | 5 | 5 | 0:15 |
| 4dfr | 229 | 0.558 | 0.0 | 0.558 | −29.4 | 1 | 0.480 | −28.3 | 4 | 5 | 0:09 |
| 1stp | 601 | 1.054 | 0.0 | 1.054 | −25.6 | 1 | 0.578 | −24.2 | 27 | 5 | 0:15 |
| 2gbp | 197 | 0.406 | 0.0 | 0.406 | −43.5 | 1 | 0.116 | −38.5 | 5 | 5 | 0:06 |
| 1abe | 601 | 0.474 | 0.0 | 0.474 | −35.8 | 1 | 0.571 | −33.2 | 9 | 5 | 0:13 |
| 1dbb | 74 | 6.762 | 8.2 | 6.739 | −13.6 | 3 | 0.467 | −27.7 | 1 | 0 | 0:13 |
| 3tpi | 206 | 1.006 | 0.0 | 1.006 | −49.6 | 1 | 0.743 | −46.9 | 26 | 5 | 0:21 |
| 6rsa | 441 | 1.365 | 0.0 | 1.365 | −43.3 | 1 | 0.693 | −45.9 | 1 | 5 | 0:33 |
| 1ack | 578 | 0.448 | 0.0 | 0.448 | −30.3 | 1 | 0.374 | −27.8 | 38 | 4 | 0:09 |
| 1tdb | 318 | 1.757 | 0.0 | 1.757 | −28.1 | 1 | 1.400 | −24.0 | 5 | 5 | 0:19 |
| 1acj | 168 | 3.521 | 1.3 | 0.862 | −24.6 | 3 | 0.585 | −24.9 | 3 | 3 | 0:07 |
| 2ak3 | 256 | 6.780 | 7.9 | 3.246 | −26.7 | 37 | 0.781 | −34.6 | 1 | 0 | 0:28 |
| 1baf | 236 | 4.659 | 3.8 | 2.316 | −21.1 | 6 | 0.378 | −22.9 | 4 | 0 | 0:11 |
| 2mth | 466 | 13.142 | 5.4 | 1.775 | −11.4 | 90 | 0.867 | −12.6 | 59 | 0 | 0:09 |
| 1mup | 301 | 4.456 | 1.5 | 0.810 | −14.5 | 8 | 0.393 | −13.9 | 41 | 0 | 0:08 |
| 4fab | 25 | 4.774 | 7.1 | 2.622 | −17.5 | 4 | 1.258 | −20.4 | 3 | 0 | 0:03 |
| 1hdc | 211 | 8.863 | 5.4 | 1.929 | −26.5 | 13 | 0.234 | −25.9 | 15 | 0 | 0:12 |

[a] All energies are in kcal/mol.
[b] Difference in energy between the lowest energy solution and the best prediction.
[c] Elapsed time in minutes and seconds on a Linux workstation (600 MHz).

tween the benzamidine $NH_2^+$ groups and Asp 189 Oδ1 and Oδ2 are shortened at the expense of the main chain hydrogen bond to Gly 219 O. Also the orientation of the benzene ring of the substrate is altered slightly (see Figure 4). This is a consistent difference between the experimental structure and all the top five predicted conformations using either algorithm. Shortening of the hydrogen bonds between the ligand and Asp 189 is also observed in the minimized crystal conformation, suggesting that the force field has an influence on this predicted difference. The absence of binding site waters HOH 416 (which hydrogen bonds to one of the ligand $NH_2^+$ groups) and HOH 710 (which is in van der Waals contact with the benzene ring) could account for the small differences observed.

*Analysis of predictions in dataset 2*

The complexes were chosen because they have previously been identified as problematic [3]. As expected the level of success is lower and the results for individual complexes vary more than those of dataset 1. Three of the protein-ligand complexes 6rsa, 1ack and 1tdb are reproduced in excellent agreement with the observed conformations with both algorithms. Previously failure for 6rsa and 1ack was attributed to an underestimation of the hydrophobic effect in binding. For 1tdb poor geometry of the ligand was regarded as being a contributing factor (although the ligand phosphate group has not been modeled in this study). There are four successes from the ten complexes using the pose clustering algorithm, but only three using the geometric hashing method. 1acj is also a convincing success with the pose clustering algorithm. The failure in the remaining complexes can be separated primarily into two groups in which:

(i) Low energy/RMSD solutions are not generated by the docking algorithm, however, the minimized experimental structure is the top solution (or close) i.e. if the crystal conformation had been generated by the algorithm it would have correctly predicted the binding mode. This represents failure in the complexes 1acj (with the geometric hashing algorithm), 2ak3 and to a lesser extent the complexes 1baf and 4fab.

(ii) Low energy/RMSD solutions are not generated, also, the minimized experimental structure

*Table 2.* Docking of Datasets 1 and 2 using the pose clustering algorithm.

| PDB complex | No. of solutions | Lowest energy[a] | | Best prediction | | | Minimized experimental | | | No. < 2 Å in Top 5 | Elapsed time[c] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | RMSD | $\Delta\Delta E_{LE \to BP}$[b] | RMSD | $\Delta E$ | Rank | RMSD | $\Delta E$ | Rank | | |
| 1ldm | 355 | 0.699 | 0.0 | 0.699 | −36.6 | 1 | 0.424 | −34.1 | 3 | 5 | 0:06 |
| 3ptb | 280 | 0.896 | 0.0 | 0.896 | −37.9 | 1 | 0.347 | −32.5 | 101 | 5 | 0:06 |
| 2phh | 638 | 0.518 | 0.0 | 0.518 | −34.3 | 1 | 0.375 | −30.0 | 26 | 5 | 0:11 |
| 1ulb | 657 | 0.664 | 0.0 | 0.664 | −28.3 | 1 | 0.147 | −28.0 | 3 | 5 | 0:13 |
| 4dfr | 325 | 0.942 | 0.0 | 0.942 | −28.8 | 1 | 0.480 | −28.3 | 3 | 3 | 0:09 |
| 1stp | 797 | 0.739 | 0.0 | 0.739 | −25.7 | 1 | 0.578 | −24.2 | 23 | 5 | 0:15 |
| 2gbp | 301 | 0.475 | 0.0 | 0.475 | −44.0 | 1 | 0.116 | −38.5 | 3 | 5 | 0:07 |
| 1abe | 401 | 0.744 | 0.0 | 0.744 | −35.0 | 1 | 0.571 | −33.2 | 5 | 5 | 0:08 |
| 1dbb | 333 | 1.070 | 0.0 | 1.070 | −26.6 | 1 | 0.467 | −27.7 | 1 | 2 | 0:09 |
| 3tpi | 180 | 1.230 | 0.0 | 1.230 | −49.7 | 1 | 0.743 | −46.9 | 17 | 5 | 0:09 |
| 6rsa | 458 | 1.192 | 0.0 | 1.192 | −43.1 | 1 | 0.693 | −45.9 | 1 | 4 | 0:17 |
| 1ack | 702 | 0.649 | 0.0 | 0.649 | −30.9 | 1 | 0.374 | −27.8 | 63 | 3 | 0:11 |
| 1tdb | 507 | 1.534 | 0.0 | 1.534 | −26.2 | 1 | 1.400 | −24.0 | 4 | 5 | 0:14 |
| 1acj | 317 | 0.663 | 0.0 | 0.663 | −24.9 | 1 | 0.585 | −24.9 | 3 | 4 | 0:10 |
| 2ak3 | 325 | 7.524 | 5.9 | 2.667 | −25.0 | 53 | 0.781 | −34.6 | 1 | 0 | 0:13 |
| 1baf | 527 | 4.271 | 7.8 | 2.931 | −18.4 | 19 | 0.378 | −22.9 | 8 | 0 | 0:12 |
| 2mth | 693 | 4.813 | 4.3 | 1.209 | −11.9 | 114 | 0.867 | −12.6 | 88 | 0 | 0:12 |
| 1mup | 350 | 3.026 | 1.5 | 0.734 | −13.7 | 45 | 0.393 | −13.9 | 39 | 0 | 0:07 |
| 4fab | 165 | 5.275 | 2.9 | 2.635 | −25.1 | 2 | 1.258 | −20.4 | 29 | 0 | 0:06 |
| 1hdc | 572 | 9.029 | 5.2 | 1.498 | −26.4 | 18 | 0.234 | −25.9 | 25 | 0 | 0:19 |

[a,b,c] See legend to Table 1.

scores poorly i.e. if the crystal conformation had been generated by the algorithm it would still not have correctly predicted the binding mode. This appears to be the case in 2mth, 1mup and 1hdc.

Those cases in group (i) would appear to be a failure of the algorithm to find the energy minimum occupied by the experimental structure, whilst (ii) would appear to be a failure of the energy function to recognise the experimental structure as the global minimum. This may be a result of the failure of the energy function to reproduce the energetics of the system i.e. neglect of explicit solvent or inadequate treatment of the hydrophobic effect as originally suggested for failures in 6rsa, 1ack, 1acj and 4fab or in the case of 1mup due to the absence of any hydrogen bonding [3]. Alternatively, as previously suggested for 2mth low resolution of the experimental structure may give rise to poor geometry which may in turn cause the calculated interaction energies to be unstable.

In order to test these hypotheses further investigation of the dependence of the results on the docking parameters has been carried out for complexes in dataset 2. In development of the algorithm the parameters that primarily influence success are (1) the cut-off distance, $R_{cut}$, between minima used in defining interaction sites in both algorithms and (2) the value of the maximum differences in distance between matched ligand atom-atom and interaction site grid-grid distances given by δ (used in the pose clustering algorithm) or *binsize* (used in the geometric hashing algorithm). In either case the values can be systematically varied to find optimal parameters for any particular protein-ligand complex. Table 3 shows the results of this type of analysis using the pose clustering algorithm (using slightly different energy minimization parameters from the solutions given in Tables 1 and 2).

The values of $R_{cut}$ and δ are systematically varied between 1.0 and 3.0 in steps of 0.5 (i.e. there are 25 different combinations for each protein-ligand complex). The results are then pooled to create a combined ranking. It can be seen that both 2ak3 and 1baf now have top ranked solutions that are close to the experimental structure. Not only this but they also have between three and all five close solutions in the top five (and both have thirteen close solutions in the top twenty), which represents a convincing success for the placement algorithm when compared with re-

*Table 3.* Docking of Dataset 2 using the pose clustering algorithm.

| PDB complex | Lowest energy[a] | | Best prediction | | | Minimized experimental | | | No. < 2 Å Top 5[c] & Top (20)[c] |
|---|---|---|---|---|---|---|---|---|---|
| | RMSD | $\Delta\Delta E_{LE \to BP}$[b] | RMSD | $\Delta E$ | Rank[c] | RMSD | $\Delta E$ | Rank[c] | |
| 6rsa | 1.219 | 0.0 | 1.219 | −45.2 | 1 | 1.655 | −38.7 | 29 | 5 (19) |
| 1ack | 3.088 | 0.1 | 0.904 | −29.4 | 2 | 0.804 | −28.3 | 86 | 2 (12) |
| 1tdb | 1.729 | 0.0 | 1.729 | −26.9 | 1 | 0.370 | −17.7 | 112 | 5 (19) |
| 1acj | 2.667 | 2.4 | 0.533 | −24.9 | 15 | 0.547 | −24.5 | 21 | 0 (1) |
| 2ak3 | 1.001 | 0.0 | 1.001 | −34.3 | 1 | 1.051 | −32.6 | 5 | 5 (13) |
| 1baf | 1.236 | 0.0 | 1.236 | −26.1 | 1 | 1.080 | −22.2 | 131 | 3 (13) |
| 2mth | 9.405 | 2.9 | 1.415 | −12.3 | 161 | 1.281 | −11.8 | 236 | 0 (0) |
| 1mup | 4.076 | 1.3 | 1.614 | −14.5 | 43 | 0.995 | −14.1 | 118 | 0 (0) |
| 4fab | 3.122 | 8.0 | 1.762 | −21.5 | 463 | 0.950 | −18.6 | 823 | 0 (0) |
| 1hdc | 13.436 | 6.5 | 2.530 | −24.7 | 21 | 0.409 | −24.2 | 25 | 0 (0) |

[a,b]See legend to Table 1.
[c]Rank takes into account degeneracy.

sults in Table 2. Success is also maintained with the complexes 6rsa, 1ack and 1tdb. The docking of the other five protein-ligand complexes become more convincing failures when judged by the ranking of their best predicted structure. The lower ranking of the minimized experimental structures of 1acj and 4fab (classified as group (i) above) than those of generated low energy/RMSD solutions now indicates group (ii) status. This later study demonstrates that although considerable success is achieved for dataset 1 using the optimized set of search parameters ($R_{cut}$, δ and *binsize*) the optimal search parameters are in fact receptor/ligand dependent. This is clearly the case for 2ak3 and 1baf. A useful analogy might be to think of the receptor-ligand interaction as a pharmacophore model in which critical features (i.e. distances and probe interaction types) of the system are the important factors in determining the success of the model. Clearly a method that pre-optimizes the parameters using prior knowledge of the receptor and ligand without the need to systematically test all combinations would further enhance the search strategy employed here. This will be the subject of further study.

*Comparison of the geometric hashing and pose clustering methods*

Comparison of Tables 1 and 2 indicates that the pose clustering algorithm is slightly superior to the geometric hashing algorithm. The number of successful dockings is greater both in dataset 1 and dataset 2 (as determined by the best prediction and the number of solutions < 2 Å ranked in the top five). The pose

clustering algorithm generates more solutions with the same constraints (geometric, bumps, degeneracy and RMSD difference < 0.75 Å) with on average ∼ 60% more solutions in dataset 1 and 100% in dataset 2 (with a large standard deviation). Certain ligands (e.g. 1dbb, 4fab) are problematic when using the geometric hashing algorithm with very few solutions generated. These ligands are large with relatively restricted binding sites therefore relatively few solutions pass the bumps test.

The pose clustering placement algorithm is faster in all cases (results not shown). Several larger ligands have much longer run times with the geometric hashing algorithm (e.g. compare times for 3tpi, 6rsa and 2ak3). This is because unlike the pose clustering method that generates receptor triplets on-the-fly the geometric hashing method creates a hash table based on all triplet combinations prior to matching. For large ligands with many probe types this can be a long list. The ligands 3tpi, 6rsa and 2ak3 are three such cases where building the receptor hash table takes most of the CPU time at between 13–19 s whilst the pose clustering takes 2–3 s for these ligands. The time differential is much less distinct for smaller ligands with few probe types. Invariably with the pose clustering method the rate-limiting step is energy minimisation of generated solutions. Often in the literature the hash table generation time is not quoted in the total run time. If we neglect this generation time then geometric hashing is faster than pose clustering.

**Discussion**

A new computer algorithm Q-fit is described here that preferentially samples low energy modes of ligand binding using a probabilistic scheme in a fully automated way. The method uses a presorted list of point interaction energies to choose optimal atom triplets on the ligand that match to receptor triplets in the receptor binding site such that the interaction energy is maximised. The two search methods proved to be very successful in their ability to predict the experimental binding modes when using the point interaction model and force field of GRID [16] and an energy minimization method. The binding modes of several ligand fragments are predicted in dataset 2 that have previously proved problematic. The binding modes of ligands used to bench-mark several other docking algorithms (dataset 1) are reproduced in excellent agreement with experiment. The results compare favourably with the results of other docking algorithms in the literature where RMSD is used as the criterion for success. The number of close solutions in the top five ranked solutions is generally high in cases where the ligand is docked successfully. This gives an indication of the reliability of the method in terms of the reproducibility of finding low energy/low RMSD solutions that have been generated independently of one another. However, it must be stressed that direct comparisons of this nature are difficult due to differences in parameterization, inclusion of ligand flexibility, the size of the ligand fragments docked and extent of the receptor binding site searched etc. Docking of a single ligand fragment including energy minimization can be performed in 10–30 s on an standard Linux workstation using conservative parameters for the search depth. Whilst efforts have concentrated on looking at the ability of the method to reproduce experimental binding locations i.e. 'the docking problem' a high level of success can still be achieved with much lower search depths and at less computational cost (results not shown). In most cases the minimization step is the most computationally demanding part of the process therefore CPU is approximately linearly related to the search depth (assuming that the same proportion of placements pass the bumps test). Improvements in the efficiency of the minimization algorithm are an obvious area for further investigation. A particularly important aspect of the scheme used here is the use of the single point probe model of GRID. This allows heavy atom placements at particular discrete points in the receptor binding site, since it is assumed that the probe types can optimally orientate themselves to satisfy hydrogen bonding requirements with the receptor. However, this may not always be the case, as the hydrogen bond geometry of the ligand may not be optimal in the generated complex. It should also be noted that whilst discretization might influence the results of this and other grid based methods [9, 10, 23, 27, 36] the effects are generally considered to be too small to play a significant role at small grid spacings (used here and in the other programs). The high level of success obtained here shows the current probe model works well.

Implicit in the model is the assumption of rigid receptor and maintenance of the small molecule geometry. The method has only been tested on rigid fragments or those with few rotatable bonds. The aim is to provide a fast placement algorithm for small fragments. This could be used in the combinatorial design of new ligands. Application to larger multi-component ligands will require a flexible treatment of the ligand. Indeed we see this current implementation as the first step in developing a flexible protein-ligand docking method. The results obtained here are very encouraging, however, the scoring function does not contain terms for ligand or protein desolvation on binding except those included implicitly in the solvent screened interaction energy. This is expected to underestimate the importance of the hydrophobic effect and overestimate the importance or charge-charge interactions in estimating binding energy. Future developments might include a more sophisticated solvation scheme such as continuum electrostatic models. Methods such as the combined Poisson-Boltzmann/surface area models [31, 32] are probably be too slow for this purpose. However, the Generalised Born approximation [33] has recently been implemented in molecular docking methods [34, 35]. The user must also make decisions about the protonation state of both the protein and ligand and the approximate location of the ligand binding site. In addition receptor bound water molecules are ignored here. Their inclusion is clearly important if prior experimental information exists, and could lead to better predicted binding geometries in some systems. Therefore no claims are made as to the effectiveness of this approach to calculating binding free energy.

## Acknowledgements

## References

1. Kuntz, I.D., Blaney, J.M., Oatley, S.J., Langridge, R. and Ferrin, T.E., J. Mol. Biol., 161 (1982) 269.
2. Rarey, M., Kramer, B., Lengauer, T. and Klebe, G., J. Mol. Biol., 261 (1996) 470.
3. Jones, G., Willett, P., Glen, R.C., Leach, A.R. and Taylor, R., J. Mol. Biol., 267 (1997) 727.
4. Miranker, A. and Karplus, M., Proteins, 11 (1991) 29.
5. Bohm, H.J., J. Comput. Aid. Mol. Des., 6 (1992) 61.
6. Gillet, V., Johnson, A.P., Mata, P., Sike, S. and Williams, P., J. Comput. Aid. Mol. Des., 7 (1993) 127.
7. Eisen, M.B., Wiley, D.C., Karplus, M. and Hubbard, R.E.,. Proteins, 19 (1994) 199.
8. Knegtel, R.M. and Wagener, M., Proteins, 37 (1999) 334.
9. Goodsell, D.S. and Olson, A., Proteins, 8 (1990) 195.
10. Hart, T.N. and Read, R.J., Proteins, 13 (1992) 206.
11. Judson, R.S., Jaeger, E.P. and Treasurywala, A.M., J. Mol. Struct., 308 (1994) 191.
12. 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.
13. Baxter, C.A., Murray, C.W., Clark, D.E., Westhead, D.R. and Eldridge, M.D. Proteins, 33 (1998) 367.
14. Lengauer, T. and Rarey, M., Curr. Opin. Struct. Biol., 6 (1996) 402.
15. Dixon, J.S., Proteins, S1 (1997) 198.
16. Goodford, P.J., J. Med. Chem., 28, (1985) 849.
17. Moon, J.B. and Howe, W.J., Proteins, 11 (1991) 314.
18. Marrone, T.J., Briggs, J.M. and McCammon, J.A., Annu. Rev. Pharmacol. Toxicol., 37 (1997) 71.
19. Jackson, R.M., Gabb, H.A. and Sternberg, M.J., J. Mol. Biol., 276 (1998) 265.
20. Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N. and Bourne, P.E., Nucleic Acids Res., 28 (2000) 235.
21. Bachar, O., Fischer, D., Nussinov, R. and Wolfson, H., Protein Eng., 6 (1993) 279.
22. Olson, C.F. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. IEEE Computer Society Press, Seattle, Washington, 1994, pp. 251–258.
23. Lawrence, M.C. and Davis, P.C., Proteins, 12 (1992) 31.
24. McLachlan, A.D., J. Mol. Biol., 128 (1979) 49.
25. Schnecke, V. and Kuhn, L.A., Intell. Syst. Mol. Biol., 7 (1999) 242.
26. Nelder, J.A. and Mead, R., Comput. J., 7 (1965) 308.
27. Gschwend, D.A. and Kuntz, I.D., J. Comput. Aid. Mol. Des., 10 (1996) 123.
28. Boobbyer, D.N., Goodford, P.J., McWhinnie, P.M. and Wade, R.C., J. Med. Chem., 32 (1989) 1083.
29. Morris, G.M., Goodsell, D.S., Halliday, R.S., Huey, R., Hart, W.E., Belew, R.K. and Olson, A.J., J. Comput. Chem., 19 (1998) 1639.
30. Marquart, M., Walter, J., Deisenhofer, J., Bode, W. and Huber, R., Acta Crystallogr. Sect. B., 39 (1983) 480.
31. Gilson, M.K. and Honig, B., Proteins, 4 (1988) 7.
32. Jackson, R.M. and Sternberg, M.J., J. Mol. Biol., 250 (1995) 258.
33. Still, W.C., Tempczyk, A., Hawley, R.C. and Hendrickson, T., J. Am. Chem. Soc., 112 (1990) 6127.
34. Given, J.A. and Gilson, M.K., Proteins, 33 (1998) 475.
35. Majeux, N., Scarsi, M., Apostolakis, J., Ehrhardt, C. and Caflisch, A., Proteins, 37 (1999) 88.
36. Meng, E.C, Shoichet, B.K. and Kuntz, I.D., J. Comput. Chem. 13 (1992) 505.
37. Boehm, H-J., Boehringer, M., Bur, D., Gmuender, H., Huber, W., Klaus, W., Kostrewa, D., Kuehne, H., Luebbers, T., Nathalie Meunier-Keller, N. and Mueller, F., J. Med. Chem., 43 (2000), 2664.