



## DOCKING METHODS

*Abhilash M.*

Department of Biotechnology, The Oxford college of Engineering, Bangalore, INDIA

Corresponding author abhibiotek@gmail.com

### **ABSTRACT**

The binding of small molecule ligands to large protein targets is central to numerous biological processes. The accurate prediction of the binding modes between the ligand and protein, (the docking problem) is of fundamental importance in modern structure-based drug design. An overview of current docking techniques is presented with a description of applications including single docking experiments and the virtual screening of databases.

### **KEY WORDS:**

Lead optimisation, molecular docking, protein-ligand complexes, scoring, virtual screening

### **INTRODUCTION**

The number of algorithms available to assess and rationalise ligand protein interactions is large and ever increasing. Many algorithms share common methodologies with novel extensions, and the diversity in both their complexity and computational speed provides a plethora of techniques to tackle modern structure based drug design problems<sup>1</sup>. Assuming the receptor structure is available, a primary challenge in lead discovery and optimisation is to predict both ligand orientation and binding affinity; the former is often referred to as 'molecular docking'<sup>2</sup>. The algorithms that address this problem have received much attention<sup>3</sup>, indicating the importance of docking to a drug design project. Owing to the increase in computer power and algorithm performance, it is now possible to dock thousands of ligands in a timeline which is useful to the pharmaceutical industry<sup>4</sup>

a drug design project. Owing to the increase in computer power and algorithm performance, it is now possible to dock thousands of ligands in a timeline which is useful to the pharmaceutical industry<sup>4</sup>.

Despite the large size of this field, we have attempted to summarise and classify the most important docking methods. The principal techniques currently available are: molecular dynamics, Monte Carlo methods, genetic algorithms, fragment-based methods, point complementarity methods, distance geometry methods, tabu searches and systematic searches. Algorithm examples and the test cases used to validate the models will be discussed

#### *Large scale docking and virtual screening*

Molecular docking is often used in virtual screening methods<sup>5</sup>, whereby large virtual libraries of compounds are reduced in size to a manageable subset, which, if successful, includes molecules



## DOCKING METHODS

with high binding affinities to a target receptor. The potential for a docking algorithm to be used as a virtual screening tool is based on both speed and accuracy. This review will therefore highlight those docking methods that have been used in virtual screening applications.

### *Docking and de novo design methods*

For the purpose of this review, a broad distinction is made between docking algorithms and *de novo* design methods. This is arguably subjective and in many cases significant overlap in methodology occurs between the two strategies. Examples of *denovo* design tools are BUILDER<sup>6</sup>, CONCEPTS<sup>7</sup>, CONCERTS<sup>8</sup>, DLD/MCSS<sup>9</sup>, Genstar<sup>10</sup>, Group-Build<sup>11</sup>, Grow<sup>12</sup>, HOOK<sup>13</sup>, Legend<sup>14</sup>, LUDI<sup>15</sup>, MCDNLG<sup>16</sup>, SMOG<sup>17</sup> and SPROUT<sup>18</sup>. LUDI is given as an example of a *denovo* design tool applied to the docking problem.

### *Search algorithms*

A rigorous search algorithm would exhaustively elucidate all possible binding modes between the ligand and receptor. All six degrees of translational and rotational freedom of the ligand would be explored along with the internal conformational degrees of freedom of both the ligand and protein. However, this is impractical due to the size of the search space. For a simple system<sup>19</sup> comprising a ligand with four rotatable bonds and six rigid-body alignment parameters, the search space has been estimated as follows. The alignment parameters are used to position the ligand relative to the protein in a cubic active site measuring 103 Å<sup>3</sup>. If the angles are considered in 10 degree increments and translational parameters on a 0.5 Å grid there are approximately 4×10<sup>8</sup> rigid body degrees of freedom to sample, corresponding to 6×10<sup>14</sup> configurations (including the four rotatable torsions) to be searched. This would require approximately 2 000 000 years of computational time at a rate of 10 configurations

per second. As a consequence only a small amount of the total conformational space can be sampled, and so a balance must be reached between the computational expense and the amount of the search space examined.

The practical application of such an extensive search involves the sampling of many high energy unfavourable states which can restrict the success of an optimisation algorithm. In practice therefore, to sample such a large search space the computational expense is limited by applying constraints, restraints and approximations to reduce the dimensionality of the problem in an attempt to locate the global minimum as efficiently as possible. A common approximation in early docking algorithms was to treat both the ligand and target as rigid bodies and only the six degrees of translational and rotational freedom were explored. One of the first examples of such an algorithm is the early implementation of the program DOCK<sup>20</sup> (see Fragment-Based Methods). Although these methods have been successful in certain cases<sup>21</sup>, there is a limitation to the rigid body docking paradigm in that the ligand conformation must be close to the experimentally observed conformation when bound to the target. Furthermore, numerous examples of conformational change of the target upon binding, for example the binding of cyclosporin A to cyclophilin<sup>22</sup>, have led the drive to incorporate conformational flexibility into the search algorithm.

A common approach in modelling molecular flexibility is to consider only the conformational space of the ligand, assuming a rigid receptor throughout the docking protocol. The techniques used to incorporate conformational flexibility into a docking protocol will be discussed in some detail. However, the searching algorithm is only half the docking problem; the other factor to



## DOCKING METHODS

be incorporated into a docking protocol is the scoring function.

### *Scoring functions*

Generating a broad range of binding modes is ineffective without a model to rank each conformation that is both accurate and efficient. The scoring function should be able to distinguish the experimental binding modes from all other modes explored through the searching algorithm. A rigorous scoring function will generally be computationally expensive and so often the function's complexity is reduced, with a consequential loss of accuracy. Scoring methods can range from molecular mechanics force fields such as AMBER<sup>23</sup>, OPLS<sup>24</sup> or CHARMM<sup>25</sup>, through to empirical free energy scoring functions<sup>26</sup> or knowledge based functions<sup>27</sup>. The currently available docking methods utilise the scoring functions in one of two ways. The first approach uses the full scoring function to rank a proteinlig and conformation. The system is then modified by the search algorithm, and the same scoring function is again applied to rank the new structure. The alternativemethod is to use a two stage scoring function. In this approach a reduced function is used to direct the search strategy and a more rigorous scoring function is then used to rank the resulting structures. These directed methods make assumptions about the energy hypersurface, often omitting computationally expensive terms such as electrostatics and considering only a few types of interaction such as hydrogen bonds. Such algorithms are therefore directed to areas of importance as determined by the reduced scoring function. Examples of directed methods are GOLD<sup>28</sup> and DOCK<sup>29</sup>, and will be considered in more detail in the following sections.

A serious limitation in many existing scoring functions is the tendency to either neglect solvation effects or use solvent models in a snap-

shot fashion. A snapshot method involves the generation of structures *in vacuo*, that are subsequently ranked with a scoring function that includes a solvent model. The search function is therefore directed to the conformational space which favours the *in vacuo* conformations. Furthermore, the structural role of bound solvent molecules and ions is often not considered, yet in the HIV-1 protease<sup>30</sup> system for example, it has been shown that explicit waters play an important role in ligand binding<sup>31</sup>. A brief description of the scoring and searching function will be given for each docking method in the following sections. The core components of the algorithm will be described, with a brief synopsis of the test cases used to validate the algorithms.

### *Molecular dynamics*

There are many programs to perform molecular dynamics (MD) simulations such as AMBER<sup>32</sup> and CHARMM<sup>25</sup>. MD involves the calculation of solutions to Newton's equations of motions. Using standard MD to find the global minimum energy of a docked complex is difficult since traversing the rugged hypersurface of a biological system is problematic. Often an MD trajectory will become trapped in a local minimum and will not be able to step over high energy conformational barriers. Thus, the quality of the results from a standard MD simulation are extremely dependent on the starting conformation of the system.

This section focuses on novel MD techniques applied specifically to the docking problem to overcome the shortcomings of standard MD methodology. Flexible ligands have been docked to flexible receptors in solution using MD simulations by Mangoni *et al.*<sup>33</sup>, building upon the original work of Di Nola *et al.*<sup>34</sup>. The problems of obtaining adequate sampling are addressed by separating the centre of mass motion of the



## DOCKING METHODS

substrate from its internal and rotational motion. Separate thermal baths are then used for both types of substrate motion and receptor motion which permits local freezing of the various motion types.

Wang and Pak<sup>35</sup> have applied a newMD method to flexible ligand docking using a well jumping technique, where a scaling function is applied to the equations of motion to facilitate barrier crossing by effectively reducing the magnitude of the forces. Multicanonical molecular dynamics addresses the problem of limited conformational sampling and has been used as a technique to dock flexible ligands by Nakajima *et al.*<sup>36</sup> These methods operate on a single structure. However, it is common practice to generate a sub-ensemble of protein states, often using molecular dynamics, for use in docking studies. Such techniques have been summarised by Carlson and McCammon<sup>37</sup> where multiple protein structures are utilised rather than operating on a single flexible protein structure.

### Monte Carlo methods

Monte Carlo (MC) methods are among the most established and widely used stochastic optimisation techniques. The combination of atomistic potential energy models with stochastic search techniques has produced some of the most powerful methods for both structure optimisation and prediction. A significant advantage of the MC technique compared with gradient based methods, such as MD, is that a simple energy function can be used which does not require derivative information. Furthermore, through a judicious choice of move type, energy barriers can simply be stepped over. The gradient based methods are often efficient at local optimisation, but have difficulty navigating a rugged hyper surface. The standard MC method (more correctly, Metropolis MC<sup>38</sup>) involves applying random Cartesian moves to the

system and accepting or rejecting the move based on a Boltzmann probability.

Early implementations of AutoDock<sup>39,40</sup> used Metropolis MC simulated annealing with a grid based evaluation of the energy, based on the AMBER force field, to dock flexible ligands into the binding pocket of a rigid receptor. The algorithm was originally tested on six complexes and was able to reproduce the experimental binding modes, although the lowest energy structures did not always correspond to the crystallographic conformation. Prodock<sup>41</sup> uses a Monte Carlo minimisation technique to dock flexible ligands to a flexible binding site, using internal coordinates to represent the structures. This method differs from a standard MC procedure in that after each random move a local gradient-based minimisation is performed; the resulting structure is then accepted based on the Metropolis acceptance criteria. A grid based technique to evaluate the energy function is incorporated into the algorithm using Bezier splines<sup>42</sup>, which produces a smooth function that can be differentiated; this property is crucial to the local gradient-based minimisation. During the dock the magnitudes of the various potential energy terms are scaled to facilitate sampling.

The independent scaling allows the selective reduction of barriers that restrict sampling. The size of each random move is determined from an assessment of the curvature of the hyper surface using the second derivative of the energy function. Thus large moves are attempted in areas of small curvature and small moves are attempted in areas of large curvature. Two force fields are implemented in Prodock, namely AMBER<sup>23</sup> and ECEPP/3<sup>43</sup> along with a solvation model based on solvent exposed volume.

The MC method has been used to dock flexible ligands into a flexible binding site by Caflisch and coworkers<sup>44</sup>; this study built on



## DOCKING METHODS

previous work by Caflisch *et al.*<sup>45</sup> for docking an FKBP-Substrate complex. The first stage of the procedure places the ligand, at random, within the active site. This structure is then minimised *in vacuo* using a conjugate gradient minimiser with the CHARMM force field, allowing flexibility of the ligand and the protein. The Lennard-Jones and coulombic potentials are initially softened and gradually turned on throughout the course of the minimisation. This is repeated for 1000 seed structures. The seed structures are then ranked based on the potential energies calculated using the CHARMM force field. Solvation is included in the potential energy using a finite difference Poisson-Boltzmann (PB) term for the electrostatic contributions, calculated by UHBD<sup>46</sup>, and nonpolar contributions are approximated by a weighted solvent-accessible area (SA) term. The MC method is then applied to the 20 structures with the lowest energy. This implementation of the MC method (referred to as Monte Carlo minimisation or MCM), is similar to the method adopted in the program Prodock<sup>41</sup>. MCM performs conjugate gradient minimisation after each random move. The minimised structures are then accepted based on the Boltzmann acceptance criteria. The energy for each MCM stage is again calculated using the CHARMM force field with the PB/SA solvent model. Each random move samples not only the position and orientation of the ligand but also a set of randomly selected dihedrals in the ligand and in the protein. This technique has been applied to three test systems and all three produced lowest energy structures within 1.4 Å RMSD of the crystallographic structures. Caflisch and co-workers also report the importance of allowing the protein to relax upon binding of the ligand, to discriminate near-native from non-native structures. This is arguably one of the most ambitious docking projects to date.

Internal Coordinates Mechanics<sup>47</sup> (ICM) is a program to perform flexible protein-ligand docking and may be summarised as a MC minimisation method in internal coordinates. The algorithm initially makes a random move, which is one of three types; rigid body ligand move, torsion moves of the ligand, or torsion moves of the receptor side chain, using the biased probability methodology<sup>48</sup>. The side chain movement using this method is one of the defining features of the algorithm. The idea is to sample with a larger probability those regions of conformational space which are known *a priori*, based on previously defined rotamers<sup>49</sup>, to be highly populated. This is achieved by making a normally distributed step in the vicinity of the low energy rotamer states for the protein side chains. Having made a random move, local minimisation of the ECEPP/3<sup>43</sup> scoring function with a distance-dependent dielectric is performed using a conjugate gradient minimiser. An approximation for side chain entropy, loosely based around the statistical distributions of side chains, is then added to the minimised *in vacuo* ECEPP/3 energy. An electrostatic solvation term is then added to this energy, which is calculated using the MIMEL<sup>48</sup> approximation. This is a rapid approximation to the reaction field potential using the Born equation with a modification for many atoms. The modified ECEPP/3 energy is then used to test whether the structure is accepted or rejected, based on the Boltzmann criteria. A history mechanism has also been implemented to promote the discovery of new minima<sup>50</sup>.

ICM has been applied to protein-ligand docking in the CASP-2 experiments<sup>51</sup>. For the 8 complexes tested only one produced an RMSD of 1.8 Å with respect to the crystal structure; the remaining test cases were only able to give, at best, an RMSD of 3 Å. However, in most cases the prediction was reasonable; on average 50% of the



## DOCKING METHODS

ligand was docked correctly. MCDOCK<sup>52</sup> (version 1.0) applies a multiple stage strategy to dock a flexible ligand to a rigid receptor. The first stage of the docking places the ligand in the binding site. Random moves are then applied to the ligand to reduce the overlap of ligand and protein atoms. Metropolis MC with simulated annealing is then performed using a scoring function based on the CHARMM force field<sup>25</sup>. This is followed by a MC simulation which uses an adjustable temperature. In this method the temperature is increased if the acceptance ratio is too low, in an attempt to yield increased sampling. MCDOCK was tested using 19 complexes, taken from the FlexX<sup>53</sup> optimisation test set. The RMSD between the binding modes predicted by MCDOCK and the experimental binding modes, for the non-hydrogen atoms of the ligand, ranged from 0.25 to 1.84 Å.

MC simulated annealing was applied to the docking problem using HIV-1 protease inhibitors by Bouzida *et al.*<sup>54</sup> The AMBER force field was used with a desolvation correction based on the product of atomic charges and volume. To traverse efficiently the rugged energy hyper surface a soft-core smoothing function was used, for both the Lennard-Jones and electrostatic contributions to the potential energy. This methodology was used to dock two flexible ligands to a rigid X-ray structure of HIV-1 protease with some crystallographic waters retained as part of the rigid system. One of the docks reproduced the experimental binding mode. However, the second test case was not successful. The AMBER potential energy function was then exchanged for the piecewise linear potential function<sup>55</sup> (PLP). The PLP function is a simple model of ligand-protein interactions encompassing four terms: ligand and receptor nonbonded interaction terms (hydrogen bonds or steric clashes), internal torsion energies, and two penalty terms for leaving the active site and for internal clashes within the ligand. Using this

scoring function the binding modes for both test cases were successfully reproduced.

Further MC simulations were performed using 10 different protein conformations for HIV-1 protease. The method consisted of randomly moving the ligand and calculating the score for this move, using the PLP function, between the ligand and all 10 protein complexes. The lowest energy from the 10 ligand-protein combinations was then used in the MC acceptance criteria to yield a frequency distribution of binding modes. This study attempted to rationalise the population of binding modes arising from the conformational changes in both the ligand and protein. They concluded, for two ligands, that there was a high correlation between protein conformation and predicted binding mode for one ligand but that the other case showed only a weak correlation.

DockVision<sup>56</sup> is another MC based docking method, using a rigid ligand and rigid receptor. The first stage of this docking algorithm generates a random ligand orientation. The MC method is then applied to the system, except the energy function is replaced by a geometric score for atomic overlap. This is followed by an MC simulated annealing protocol using a simple potential energy function. The two stage docking procedure is then repeated for a large number of random ligand orientations. The ligand orientations generated by the MCdock are then clustered, based on a RMSD score. Two inhibitor complexes were used to test the protocol and in each case the binding geometry was correctly predicted. More recently this methodology has been applied to protein docking in CASP-2 experiments<sup>57</sup>, achieving the second highest success rate.

QXP<sup>58</sup> performs MC flexible ligand/rigid protein docking, and is part of the FLO96 package. The Metropolis MC method is initially performed on the isolated ligand using only random dihedral moves (up to 360°). This is followed by rigid body



## DOCKING METHODS

rotations and translations to align the ligand onto guide atoms within the active site. These guide atoms are simply atoms in van der Waals contact with the binding site atoms. Having aligned the atoms within the active site, the MC method is applied to the ligand using only rigid body rotations and translations. Conjugate-gradient minimisation is then performed on the ligand torsions followed by Metropolis MC on the ligand torsions. In this method a grid representation of the receptor is used. The scoring function uses the AMBER force field with short non-bonded cut-offs and a distance dependent dielectric. The original test set consisted of 12 ligand-protein complexes, with a maximum of 24 rotatable ligand dihedrals. The ligand was flexible and the receptor rigid, with single important water molecules retained in three of the complexes. Their results were compared with energy minimised structures; 11 ligands gave an RMSD of less than 0.76 Å. Affinity<sup>59</sup> is commercial program using Monte Carlo simulated annealing with a grid representation for the non-moving parts of the system [60] and an implicit representation of solvation effects<sup>61</sup>. Another commercial program is Glide<sup>62</sup> which uses a hierarchical filter to rapidly score hydrophobic and polar contacts, followed by Monte Carlo sampling with the ChemScore<sup>26</sup> scoring function.

### Genetic algorithms and evolutionary programming

Since their inception, genetic algorithms (GA) have increased in popularity as an optimisation tool. It should be noted that GAs (and evolution programming (EP)) require the generation of an initial population whereas conventional MC and MD require a single starting structure in their standard implementation. The essence of a GA is the evolution of a population of possible solutions via genetic operators (mutations, crossovers and migrations) to a final population,

optimising a predefined fitness function. Degrees of freedom are encoded into genes or binary strings and the collection of genes, or chromosome, is assigned a fitness based on a scoring function. The mutation operator randomly changes the value of a gene, crossover exchanges a set of genes from one parent chromosome to another, and migration moves individual genes from one sub-population to another.

GOLD<sup>28</sup> is a docking program that uses a GA search strategy and includes rotational flexibility for selected receptor hydrogens along with full ligand flexibility. Gene encoding is used to represent both rotatable dihedrals and ligand-receptor hydrogen bonds. A GA move operator is subsequently applied to parent chromosomes that are randomly chosen from the existing population with a bias towards the fittest members. The ligand-receptor hydrogen bonds are subsequently matched with a least squares fitting protocol to maximise the number of inter-molecular hydrogen bonds for each GA move. As a consequence the GA structure generation is biased towards inter-molecular hydrogen bonds. However each structure is ranked based on a more complex fitness function. The fitness (or scoring) function is the sum of a hydrogen bond term, a 4–8 inter-molecular dispersion potential and a 6–12 intra-molecular potential for the internal energy of the ligand. Each complex was run using an initial population of 500 individuals divided into five equal sub-populations, and migration of individual chromosomes between sub-populations was permitted. A single GA run used 100 000 genetic operations and 20 GA runs were performed. Finally, the solution with the highest fitness score was compared with the crystallographic binding mode.

AutoDock 3.0<sup>64</sup> uses a genetic algorithm as a global optimiser combined with energy minimisation as a local search method. In this



## DOCKING METHODS

implementation of AutoDock the ligand is flexible and the receptor is rigid and represented as a grid. The genetic algorithm uses two point crossover and mutation operators. For each new population a user determined fraction undergo a local search procedure using a random mutation operator where the step size is adjusted to give an appropriate acceptance ratio. The fitness function comprises five terms: a Lennard-Jones 12-6 dispersion/repulsion term; a directional 12-10 hydrogen bond term; a coulombic electrostatic potential; a term proportional to the number of sp<sup>3</sup> bonds in the ligand to represent unfavourable entropy of ligand binding due to the restriction of conformational degrees of freedom; and a desolvation term. This scoring function is based loosely around the AMBER force field from which protein and ligand parameters are taken. The desolvation term is an inter-molecular pairwise summation combining an empirical desolvation weight for ligand carbon atoms, and a pre-calculated volume term for the protein grid. Each of the five terms are weighted using an empirical scaling factor determined using linear regression analysis from a set of 30 protein-ligand complexes with known binding constants. The algorithm was originally tested on seven complexes, and for these test examples all lowest energy structures were within 1.14 Å RMSD of the crystal structure.

DIVALI<sup>65</sup> uses an AMBER-type potential energy function with a distance dependent dielectric and a genetic algorithm search function to dock four complexes. The receptor was modelled as a rigid entity and consequently a grid based energy evaluation of ligand protein interactions was performed to assess the fitness function. An additional masking operator is used that fixes part of the population which is associated with translational space so that subpopulations search different regions of the active site. Three out of

four of the test complexes gave an RMSD of 1.7 Å or less.

The program DARWIN<sup>66</sup> combines a GA and a local gradient minimisation strategy with the CHARMM-AA molecular mechanics force field, for flexible docking of three protein-carbohydrate complexes. Binary encoding is used to describe a starting ligand conformation and position; the potential energy is then locally minimised using a gradient method and the chromosome fitness is scored using the CHARMM-AA potential energy function. The populations are then modified by standard mutation and crossover operators while the protein is held rigid. Solvent contributions are assessed using a modified version of the program DelPhi<sup>67</sup> to yield finite difference solutions to the Poisson-Boltzmann equation. Although the search algorithm was able to optimise the energy landscape, certain structures were obtained with energies lower than the experimental binding mode. The false positives produced were thus attributed to limitations in the scoring function. Including specific explicit waters in the binding site increased the success of the program. The authors further note the dynamic nature of the complexes, and that multiple binding modes is a reasonable reflection of reality and not an artefact of the force field.

Judson *et al.*<sup>68</sup> were one of the first to report the application of a GA to the docking problem. A flexible ligand was used with interacting subpopulations and a gradient minimisation during the search. The method was tested by docking Cbz-GlyP-Leu-Leu into thermolysin and produced conformations which were close to the experimental binding mode, although in some cases the energies were lower than the crystal conformation.

Gehlhaar *et al.*<sup>55</sup> have applied evolutionary algorithms to flexible ligand docking in an HIV-1 protease complex, using the previously described





## DOCKING METHODS

PLP scoring function<sup>55</sup>. An initial population is generated and the fitness of each member is evaluated based on the scoring function. The fitness of the members are then compared with a predetermined number of opponent members chosen at random. The members are then ranked by the number of wins, and the highest ranking solutions are chosen as a new population. All surviving solutions are used to produce offspring with a mutation operator, such that the population size is constant. This protocol is repeated until a user defined number of iterations is exceeded; a conjugate gradient optimisation is then performed on the best member. Interestingly, for this test case, previous docking attempts have failed. This failure was attributed to high energy barriers<sup>69</sup>. Consequently, the repulsive term was slowly turned on through the course of the simulation, in an analogous fashion to MC simulated annealing. 100 simulations were run and the crystal structure was reproduced 34 times with a maximum RMSD of 1.5 Å; these solutions were the lowest energy docks.

### *Fragment-based methods*

The broad philosophy of fragment based docking methods can be described as dividing the ligand into separate portions or fragments, docking the fragments, followed by the linking of fragments. These methods require subjective decisions on the importance of the various functional groups in the ligand, which can result in the omission of possible solutions, due to assumptions made about the potential energy landscape. Furthermore, a judicious choice of base fragment is essential for these methods, and can significantly affect the quality of the results. The docking of fragments and the subsequent joining of the docked fragments has been widely used in *de novo* design methods. Although, for this review, only a few *de novo* programs have been

considered, there is a considerable overlap of methodologies.

One of the most popular programs to perform fragment docking is the incremental construction algorithm FlexX<sup>53</sup>. The initial phase is the selection of the base fragment for the ligand from which possible conformations are formed based on the MIMUMBA<sup>70</sup> torsion angle database. As with all fragment based methods the choice of base fragment is crucial to the algorithm; it must contain the predominant interactions with the receptor. Early implementations of FlexX required manual selection of this base fragment but this process has been subsequently automated<sup>71</sup>. Following the selection of the base fragment an alignment procedure is performed to optimise the number of favourable interactions. These interactions are based primarily on hydrogen bond geometric constraints but also include hydrophobic interactions. In this stage, the base fragment is considered rigid, and three sites on the fragment are mapped onto three sites of the receptor. All geometrically accessible receptor triangles are then clustered and the superposition of ligand triplets onto the receptor is performed using the method of Kabsch<sup>72</sup>. Overlaps are removed and energies are then calculated for the base fragments using Böhm's<sup>73</sup> function. Following this base fragment placement the ligand is built in an incremental fashion, where each new fragment is added in all possible positions and conformations.

Intra-molecular and inter-molecular overlaps are then removed and the placements are ranked, from which the best solutions are subjected to a clustering protocol. The highest rank solution from each cluster is then used in the next iteration. This process is repeated until the complete ligand is built, and the final structures are scored using the empirical scoring function.

The program DOCK (version 4.0<sup>29</sup>) can be summarised as a search for geometrically allowed



## DOCKING METHODS

ligand-binding modes using several steps that include: describing the ligand and receptor cavity as sets of spheres, matching the sphere sets, orienting the ligand, and scoring the orientation. New extensions to the protocol involve combining the bipartite graphs consisting of protein and ligand interaction sites, into a single docking graph where each node now represents a pairing of an atom with a site point. Clique detection is then implemented based on the methodology of Bron and Kerbosch<sup>76</sup>. The technique has been previously evaluated<sup>77, 78</sup> and was found to be the most efficient methodology for finding cliques which encode maximal pairs of interactions between matching sites. Having generated multiple orientations an inter-molecular score is calculated based, on the AMBER force field<sup>79</sup> where receptor terms are calculated on a grid. DOCK 4.0 includes ligand flexibility using a modified scoring function which incorporates an intramolecular score for the ligand<sup>80</sup>. In this version the docking is fragment based; a ligand anchor fragment is selected and placed in the receptor, followed by rigid body simplex minimisation. The conformations of the remaining parts of the ligand are searched by a limited backtrack method and minimised. This protocol was tested on 10 structures; 7 docked complexes reproduced the crystal structure with a maximum RMSD of 1.03 Å and the remaining 3 were within 1.88 Å. Although DOCK is included as a fragment based method, this is only one of several modes of operation. An alternative mode of operation is the docking of multiple random ligand conformations<sup>81, 82</sup>. This method generates a user-defined number of conformers as a multiple of the number of rotatable bonds in the ligand. If the total number of user-defined conformers is greater than the number of conformations possible, based on a set of dihedral rules, then a systematic search is performed. Otherwise the required number of conformers are generated by assigning random

dihedral values. These conformers are then docked independently. The search algorithms available in DOCK 4.0 have recently been reviewed by Ewing *et al.*<sup>82</sup>. Further extensions to DOCK have included incorporating protein flexibility using ensembles of protein structures<sup>83</sup> and the inclusion of a GB/SA<sup>84</sup> continuum model into the scoring function<sup>85</sup>.

## SUMMARY AND CONCLUSIONS

An extensive summary of currently available docking methods has been presented. Comparisons suggest that the best algorithm for docking is probably a hybrid of various types of algorithm encompassing novel search and scoring strategies. The most useful docking method will not only perform well, but will be easy to use and parametrise, and sufficiently adaptable such that different functionality may be selected, depending on the number of structures to be docked, the available computational resources, and the complexity of the problem. If the parameters cannot be generated quickly then although the algorithm may be computationally efficient, from a practical point of view it is limited. Conversely, a rapid scoring function may not necessarily be able to model some specific interactions. Algorithms that use the rigid receptor/flexible ligand approximation are well established and the most successful programs have achieved a success rate of between 70–80%. However, in the few examples where protein flexibility is incorporated into the docking algorithm, it is not clear whether the protein conformational states are sampled extensively. Furthermore, incorporating an ‘on-the-fly’ solvent model into a docking method is a problem which has only recently been addressed with varying degrees of success. Moreover, although current docking methods show great promise, fast and accurate discrimination between different ligands based on binding affinity, once



## DOCKING METHODS

the binding mode is generated, is still a significant problem.

### REFERENCES

1. Kuntz, I.D. *Science*, 257 (1992) 1078–1082.
2. Lybrand, T.P. *Curr. Opin. Struct. Biol.*, 5 (1995) 224–228.
3. Blaney, J.M. and Dixon, J.S. *Perspect. Drug Discov.*, 1 (1993) 301–319.
4. Abagyan, R. and Totrov, M. *Curr. Opin. Chem. Biol.*, 5 (2001) 375–382.
5. Walters, W.P., Stahl, M.T. and Murcko, M.A. *Drug Discovery Today*, 3 (1998) 160–178.
6. Roe, D.C. and Kuntz, I.D. *J. Comput. Aid. Mol. Des.*, 9 (1995) 269–282.
7. Pearlman, D.A. and Murcko, M.A. *J. Comput. Chem.*, 14 (1993) 1184–1193.
8. Pearlman, D.A. and Murcko, M.A. *J. Med. Chem.* 39 (1996) 1651–1663.
9. Stultz, C.M. and Karplus, M. *Proteins*, 40 (2000) 258–289.
10. Rotstein, S.H. and Murcko, M.A. *J. Comput. Aid. Mol. Des.*, 7 (1993) 23–43.
11. Rotstein, S.H. and Murcko, M.A. *J. Med. Chem.*, 36 (1993) 1700–1710.
12. Moon, J.B. and Howe, W.J. *Proteins*, 11 (1991) 314–328.
13. Eisen, M.B., Wiley, D.C., Karplus, M. and Hubbard, R.E. *Proteins*, 19 (1994) 199–221.
14. Nishibata, Y. and Itai, A. *J. Med. Chem.*, 36 (1993) 2921–2928.
15. Böhm, H.J. *J. Comput. Aid. Mol. Des.*, 6 (1992) 61–78.
16. Gehlhaar, D.K., Moerder, K.E., Zichi, D., Sherman, C.J., Ogden, R.C. and Freer, S.T. *J. Med. Chem.*, 38 (1995) 466–472.
17. Dewitte, R.S., Ishchenko, A.V. and Shakhnovich, E.I. *J. Am. Chem. Soc.*, 119 (1997) 4608–4617.
18. Gillet, V., Johnson, A.P., Mata, P., Sike, S. and Williams, P. *J. Comput. Aid. Mol. Des.*, 7 (1993) 127–153.
19. Welch, W., Ruppert, J. and Jain, A.N. *Chem. Biol.*, 3 (1996) 449–462.
20. Kuntz, I.D., Blaney, J.M., Oatley, S.J., Langridge, R. and Ferrin, T.E. *J. Mol. Biol.*, 161 (1982) 269–288.
21. Shoichet, B.K. and Kuntz, I.D. *Chem. Biol.*, 3 (1996): 151–156.
22. Wüthrich, K., Freyberg, V.B., Weber, C., Wider, G., Traber, R., Widmer, H. and Braun, W. *Science*, 254 (1991) 953–954.
23. Cornell, W.D., Cieplak, P., Bayly, C.I., Gould, I.R., Merz, K.M., Ferguson, D.M., Spellmeyer, D.C., Fox, T., Caldwell, J.W. and Kollman, P.A. *J. Am. Chem. Soc.*, 117 (1995) 5179–5197.
24. Jorgensen, W.L. and Tirado-Rives, J. *J. Am. Chem. Soc.*, 110 (1988) 1657–1666.
25. Brooks, B.R., Bruccoleri, R.E., Olafson, B.D., States, D.J., Swaminathan, S. and Karplus, M. *J. Comput. Chem.*, 4 (1983) 187–217.
26. Eldridge, M.D., Murray, C.W., Auton, T.R., Paolini, G.V. and Mee, R.P. *J. Comput. Aid. Mol. Des.*, 11 (1997) 425–445.
27. Muegge, I. and Martin, Y.C. *J. Med. Chem.*, 42 (1999) 791–804.
28. Jones, G., Willett, P., Glen, R.C., Leach, A.R. and Taylor, R. *J. Mol. Biol.*, 267 (1997) 727–748.
29. Ewing, T.J.A. and Kuntz, I.D. *J. Comput. Chem.*, 18 (1997) 1175–1189.
30. Miller, M., Schneider, J., Sathyanarayana, B.K., Toth, M.V., Marshall, G.R., Clawson, L., Selk, L., Kent, S.B.H. and Wlodawer, A. *Science*, 246 (1989) 1149–1152.
31. Lam, P.Y.S., Jadhav, P.K., Eyermann, C.J., Hodge, C.N., Ru, Y., Bachelier, L.T., Meek, J.L., Otto, M.J., Rayner, M.M., Wong, Y.N., Chang, C.H., Weber, P.C., Jackson, D.A.,

**DOCKING METHODS**

- Sharpe, T.R. and Erickson-Viitanen, S. *Science*, 263 (1994) 380–384.
32. Kollman, P.A. AMBER 5.0, University of California, San Francisco, 1996.
33. Mangoni, R., Roccatano, D. and Di Nola, A. *Proteins*, 35 (1999) 153–162.
34. Di Nola, A., Roccatano, D. and Berendsen, H.J.C. *Proteins*, 19 (1994) 174–182.
35. Pak, Y. and Wang, S. *J. Phys. Chem. B*, 104 (2000) 354–359.
36. Nakajima, N., Higo, J., Kidera, A. and Nakamura, H. *Chem. Phys. Lett.*, 278 (1997) 297–301.
37. Carlson, H.A. and McCammon, J.A. *Mol. Pharmacol.*, 57 (2000) 213–218.
38. Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H. and Teller, E., *J. Chem. Phys.*, 21 (1953) 1087–1092
39. Goodsell, D.S. and Olson, A.J. *Proteins*, 8 (1990) 195–202.
40. Morris, G.M., Goodsell, D.S., Huey, R. and Olson, A.J. *J. Comput. Aid. Mol. Des.*, 10 (1996) 293–304.
41. Trosset, J.Y. and Scheraga, H.A. *J. Comput. Chem.*, 20 (1999) 412–427.
42. Trosset, J.Y. and Scheraga, H.A. *J. Comput. Chem.*, 20 (1999) 244–252.
43. Némethy, G., Gibson, K.D., Palmer, K.A., Yoon, C.N., Paterlini, G., Zagari, A., Rumsey, S. and Scheraga, H.A. *J. Phys. Chem.*, 96 (1992) 6472–6484.
44. Apostolakis, J., Plückthun, A. and Caflisch, A. *J. Comput. Chem.*, 19 (1998) 21–37.
45. Caflisch, A., Fischer, S. and Karplus, M. *J. Comput. Chem.*, 18 (1997) 723–743.
46. Davis, M.E., Madura, J.D., Luty, B.A. and McCammon, J.A. *Comput. Phys. Commun.*, 62 (1991) 187–197.
47. Abagyan, R., Totrov, M. and Kuznetsov, D. *J. Comput. Chem.*, 15 (1994) 488–506.
48. Abagyan, R. and Totrov, M. *J. Mol. Biol.*, 235 (1994) 983–1002.
49. Ponder, J.W. and Richards, F.M. *J. Mol. Biol.*, 193 (1987) 775–791.
50. Argos, P. and Abagyan, R. *J. Mol. Biol.*, 225 (1992) 519–532.
51. Totrov, M. and Abagyan, R. *Proteins Suppl.* 1 (1997) 215–220.
52. Liu, M. and Wang, S. *J. Comput. Aid. Mol. Des.*, 13 (1999) 435–451.
53. Rarey, M., Kramer, B., Lengauer, T. and Klebe, G. *J. Mol. Biol.*, 261 (1996) 470–489.
54. Bouzida, D., Rejto, P.A., Arthurs, S., Colson, A.B., Freer, S.T., Gehlhaar, D.K., Larson, V., Luty, B.A., Rose, P.W. and Verkhivker, G.M. *Int. J. Quant. Chem.*, 72 (1999) 73–84.
55. 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) 317–324.
56. Hart, T.N. and Read, R.J. *Proteins*, 13 (1992) 206–222.
57. Janin, J. *Prog. Biophys. Molec. Biol.*, 64 (1995) 145–166.
58. McMartin, C. and Bohacek, R.S. *J. Comput. Aid. Mol. Des.*, 11 (1997) 333–344.
59. Accelrys Inc., San Diego, CA 92121.
60. Luty, B.A., Wasserman, Z.R., Stouten, P.F.W., Hodge, C.N., Zacharias, M. and McCammon, J.A. *J. Comput. Chem.*, 16 (1995) 454–464.
61. Stouten, P.F.W., Frommel, C., Nakamura, H. and Sander, C. *Mol. Simulat.*, 10 (1993) 97–120.
62. Schrödinger Inc., San Diego, CA 92122.
63. GOLD 1.2, CCDC, Cambridge, UK, 2001.
64. 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) 1639–1662.
65. Clark, K.P. and Ajay. *J. Comput. Chem.*, 16 (1995) 1210–1226.



## DOCKING METHODS

66. Taylor, J.S. and Burnett, R.M. *Proteins*, 41 (2000) 173–191.
67. Nicholls, A. and Honig, B. *J. Comput. Chem.*, 12 (1991) 435–445.
68. Judson, R.S., Jaeger, E.P. and Treasurywala, A.M. *Theochem.*, 114 (1994) 191–206.
69. Caflisch, A., Niederer, P. and Anliker, M. *Proteins*, 13 (1992) 223–230.
70. Klebe, G. and Mietzner, T. *J. Comput. Aid. Mol. Des.*, 8 (1994) 583–606.
71. Rarey, M., Kramer, B. and Lengauer, T. *J. Comput. Aid. Mol. Des.*, 11 (1997) 369–384.
72. Kabsch, W. *Acta Cryst.*, A32 (1976) 922–923.
73. Böhm, H.J. *J. Comput. Aid. Mol. Des.*, 8 (1994) 243–256.
74. Rarey, M., Kramer, B. and Lengauer, T. *Proteins*, 34 (1999) 17–28.
75. Claussen, H., Buning, C., Rarey, M. and Lengauer, T. *J. Mol. Biol.*, 308 (2001) 377–395.
76. Bron, C. and Kerbosch, J. *Comm. of the A.C.M.*, 16 (1973) 575–576.
77. Brint, A.T. and Willett, P. *J. Chem. Inform. Comput. Sci.*, 27 (1987) 152–158.
78. Grindley, H.M., Artymiuk, P.J., Rice, D.W. and Willett, P. *J. Mol. Biol.*, 229 (1993) 707–721.
79. Meng, E.C., Shoichet, B.K. and Kuntz, I.D. *J. Comput. Chem.*, 13 (1992) 505–524.
80. Makino, S. and Kuntz, I.D. *J. Comput. Chem.*, 18 (1997) 1812–1825.
81. Lorber, D.M. and Shoichet, B.K. *Protein Sci.*, 7 (1998) 938–950.
82. Ewing, T.J.A., Makino, S., Skillman, A.G. and Kuntz, I.D. *J. Comput. Aid. Mol. Des.*, 15 (2001) 411–428.
83. Knegtel, R.M.A., Kuntz, I.D. and Oshiro, C.M. *J. Mol. Biol.*, 266 (1997) 424–440.
84. Still, W.C., Tempczyk, A., Hawley, R.C. and Hendrickson, T. *J. Am. Chem. Soc.*, 112 (1990) 6127–6129.
85. Zou, X.Q., Sun, Y. and Kuntz, I.D. *J. Am. Chem. Soc.*, 121 (1999) 8033–8043.