

Docking Ligands into Flexible and Solvated Macromolecules. 1. Development and Validation of FITTED 1.0

Christopher R. Corbeil, Pablo Englebienne, and Nicolas Moitessier*

Department of Chemistry, McGill University, 801 Sherbrooke St W, Montreal, Quebec, Canada H3A 2K6

Received June 22, 2006

We report the development and validation of a novel suite of programs, FITTED 1.0, for the docking of flexible ligands into flexible proteins. This docking tool is unique in that it can deal with both the flexibility of macromolecules (side chains and main chains) and the presence of bridging water molecules while treating protein/ligand complexes as realistically dynamic systems. This software relies on a genetic algorithm to account for the flexibility of the two molecules as well as the location of bridging water molecules. In addition, FITTED 1.0 features a novel application of a switching function to retain or displace key water molecules from the protein–ligand complexes. Two independent modules, ProCESS and SMART, were developed to set up the proteins and the ligands prior to the docking stage. Validation of the accuracy of the software was achieved via the application of FITTED 1.0 to the docking of inhibitors of HIV-1 protease, thymidine kinase, trypsin, factor Xa, and MMP to their respective proteins.

INTRODUCTION

Fast, cost-effective, and accurate methods of drug design are essential to modern day medicinal chemistry. Docking-based rational design methods provide a quick and economical alternative to high throughput screening or more traditional drug discovery and are increasingly popular alternatives.¹

Over the past few years, several comparative studies of docking programs have been published and show the poor accuracy of some of the commercially available packages, with Glide and GOLD being among the best programs.^{2–8} In most studies, inhibitors are accurately docked back to their corresponding protein structure (self-docking). However, it has been shown that docking to other structures (cross-docking) performs poorly.^{9–11} This failure results in part from the use of inaccurate protein models. Several docking programs treat the proteins as rigid objects and do not account for conformational changes upon binding, resulting in this observed poor performance in the cross-docking studies and low enrichment factors in virtual screening studies.^{12,13} Improvement of the developed software is necessary to include more accurate protein models.

To account for the discrepancy between self- and cross-docking, various strategies have been explored and implemented in existing software.^{14,15} The program FlexE uses a set of protein structures as an input and describes the side-chain and main-chain flexibility.¹⁶ SLIDE, which handles flexible side chains,¹⁵ can also explore the main-chain flexibility when coupled with ROCK.¹⁷ Another docking program, AutoDock, models rigid proteins using grids that can be combined to approximate ensembles of conformations.^{10,18} In a fourth strategy, Glide, when merged with Prime, accounts for protein adjustments through the use of homology models.¹⁹

We have recently proposed a novel concept for the docking of ligands to solvated biopolymers,²⁰ a pharmacophore-oriented docking approach,²¹ and a genetic algorithm (GA) based docking method.²² The latter takes advantage of more than one structure to dock compounds in virtually flexible proteins. Using a similar approach to Lengauer and co-workers¹⁶ and Shoichet and co-workers,²³ we used a library of experimentally observed protein conformations and made composite structures to model the protein flexibility and to explore a wide region of conformational space. The proteins and ligands were described as genes and a mixed Lamarckian/Darwinian evolution optimized the entire complex. This virtual flexibility was found to significantly increase the accuracy of the docking of BACE-1 inhibitors.²²

We report herein the development of FITTED 1.0 (Flexibility Induced Through Targeted Evolutionary Description), a suite of programs based on a genetic algorithm (GA) with an emphasis on speed. This docking software is unique in that it can deal with both the flexibility of macromolecules and the presence of bridging “displaceable” water molecules. Additional operators to the more traditional crossover and mutations were implemented and led to a significant increase in speed. These operators simulate the learning (through energy minimization at various stages) and the early selection of individuals based on a crude estimation of their fitness (e.g., is the ligand in the binding site?). A new potential energy function modeling the interaction with displaceable water molecules and two modules (ProCESS and SMART) needed to prepare the ligands and proteins is also described. A validation of the accuracy of the docking program was performed on five different sets of protein–ligand complexes: HIV-1 protease, thymidine kinase, trypsin, factor Xa, and stromelysin-1 cocrystallized with a variety of inhibitors.

THEORY AND IMPLEMENTATION

Proof of Concept. Our previous report²² of the use of Lamarckian GA to account for both ligand and protein

*Corresponding author phone: (514)398-8543; fax: (514)398-3797; e-mail: nicolas.moitessier@mcgill.ca.

flexibility was based on Discover 3.0²⁴ as a force field engine and considered only the ligand torsion angles as degrees of freedom. The flexibility of the side chains and main chains of the target protein were modeled using a library of conformations (from available data) that could evolve by means of genetic operators (crossover and mutations). An anchor atom was needed by this early version to ensure convergence in a reasonable period of time. In practice, runs were performed in as long as 20 h for the most flexible ligands. Upon further investigations, we found that more than 96% of CPU time was consumed by intermediate minimization steps (part of the Lamarckian GA). The inclusion of ligand translational and rotational degrees of freedom led to intractable computations. This proof-of-concept led us to develop a program based on the same concept with a strong focus on the CPU time required, instead of using many independent programs that do not communicate quickly nor easily with each other. FITTED 1.0 includes a force field engine to perform conjugate gradient minimization²⁵ and a genetic algorithm.

Program Requirements and Setup. FITTED is being designed to be a docking-based virtual screening (VS) tool. Before libraries can be screened, the docking algorithm must be validated. In practice, aspects of the docking routine that are common to all runs should be performed only once. First, since protein structures are common to all runs in a VS study, it is best to have a separate program to do their setup once, quickly, and efficiently. Second, a virtual library of druglike molecules is, in practice, applied to more than one biologically relevant target and should also be prepared independently from the VS run. These two aspects led us to create modules for FITTED, namely ProCESS and SMART, described in greater detail in the following sections. The use of modules is a common practice as exemplified by the AutoDock suite of programs.²⁶ FITTED, SMART, and ProCESS can either be run from the command line in Linux or as console applications in Windows. However, the accuracy of FITTED was found to be highly compiler-dependent, and caution should be taken to ensure the suitability of the compiler before using FITTED (gcc v.3.2.3 was found to be the best).

ProCESS, a Tool To Prepare Protein Files. In order to have protein files useable by FITTED, we developed the ProCESS module (**P**rotein **C**onformational **E**nsemble **S**ystem **S**etup) which assigns the advanced residue names, advanced hydrogen names, atom types, and charges for the protein as discussed below. FITTED may use several protein files to consider the protein flexibly. However, these files must be homogeneously prepared (i.e., same atom name and number, same primary sequence). ProCESS requests all-atom proteins in mol2 format as inputs. Various programs (InsightII, Maestro, Sybyl) can be used to add hydrogen atoms to the PDB files. However, all these graphical interfaces do not generate the same mol2 files from the same input structures (various residue names, hydrogen names, order of atoms). ProCESS first ensures that the protein files are consistent and can be used unambiguously. Rules exist for the naming of atoms/groups in proteins (PDB).²⁷ For instance, CYS and ASP designate cysteine and aspartic acid residues, respectively. However, this naming does not give information on the protonation state of the side chain nor does it identify the terminal residues. CYS can be involved in a disulfide

bridge or not, while ASP can be protonated or not. As a result, these residue names cannot be used unambiguously to assign partial charges and atom types. To address this issue, the graphical interfaces assign various residue names for ionized and neutral aspartyl residues: Maestro (Schrödinger), ASP and ASH; InsightII (Accelrys), ASP- and ASP; Sybyl (Tripos), ASP and ASZ. Additional names are used for capped terminal residues in Sybyl (AMN or AMI) and CXL or CXC for the ionized or neutral terminal amino and carboxyl groups, respectively. In order to use mol2 files generated with these different interfaces, ProCESS reassigns advanced names to the added hydrogens (i.e., HA for alpha hydrogens, HB1 and HB2 for beta hydrogens) by examining their chemical environment and proceeds with the advanced residue names starting with the terminal residues. If the residue is an *N*-terminus, it is assigned a fourth letter, an N; if it is a *C*-terminus, a C is appended to the name. ProCESS then checks for the protonation state of CYS, ASP, GLU, and HIS. CYS is renamed to CYSH if not involved in a disulfide bridge; ASP and GLU are changed to ASPH and GLUH if neutral. HIS can have one of three possible names: HISE if the hydrogen is on the ϵ nitrogen, HISD if the hydrogen is on the δ nitrogen, and HISP if positively charged. By defining the chemical environment, these advanced names aid ProCESS in assigning appropriate partial charges and atom types to the protein.

In some cases, the atom ordering also varies from one PDB or mol2 file to another. To address this issue, ProCESS sorts each of the protein files by residue and atom names and checks for sequence identity; if discrepancies in the primary sequence are found, then ProCESS exits, and manual editing of the protein files is required.

As soon as the protein files are checked and made consistent, ProCESS truncates the protein input structures, using a user-defined cutoff distance around a user-defined binding site (for a list of residues, see the Supporting Information). The truncated proteins are represented as united atoms, and AMBER atom types and partial charges are assigned to each of them. The use of truncated united-atom protein structures significantly reduces the time required by FITTED to set the lists of nonbond interactions at the outset of each minimization stage. While potential energies computed using a force field are used all through the docking process, scoring of the final poses is performed by the previously developed RankScore²² scoring function. This function implicitly accounts for the entropy cost of freezing flexible residues upon binding, computed by scaling down the interactions with flexible residues as discussed below. To account for these scaling factors, ProCESS assigns new "scaled" atom types and charges derived from AMBER atom types. When processed, the protein structures are outputted in mol2 format. Similar to FITTED, ProCESS requires a keyword file which contains all the necessary parameters for ProCESS to work. A typical keyword file is given as Supporting Information.

ProCESS, a Tool To Create Binding Site Cavity Files.

As discussed below, FITTED disregards poses that are not within the binding site cavity, which is approximated by a series of overlapping spheres. The required CPU time of initial docking attempts using grids as cavity representations was not satisfactory. We have found that moving to spheres was much less CPU time-consuming (by a factor of 3).

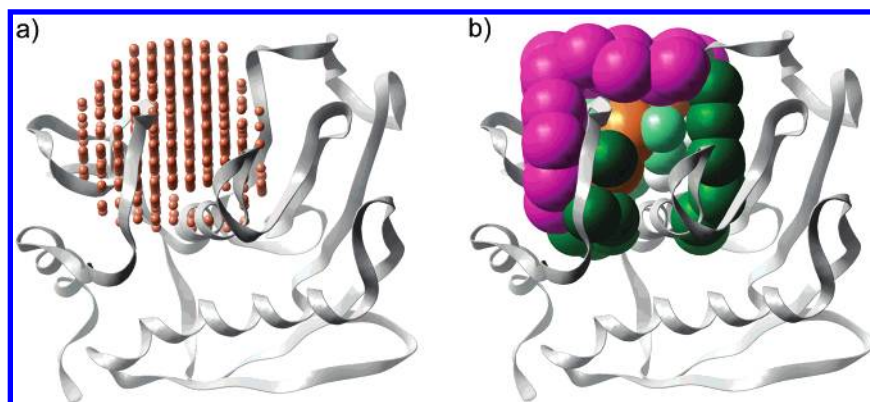


Figure 1. The binding site of 1d8m mapped as (a) a set of points and (b) a set of spheres, the spheres are colored by size range (from 1.5 Å to over 4.0 Å).

ProCESS creates the overlapping spheres by first generating an evenly spaced grid and keeping the points that do not clash with the protein (Figure 1a). If more than one protein is used, then the point must clash with all proteins to be removed. Thus, alternative accessible spaces are maintained.

Next the points are converted into spheres (Figure 1b). For this purpose, each sphere is inflated until making contact with a protein atom center (slightly overlapping with the protein surface) or one of the grid edges. The obtained sphere size and center are archived. Smaller spheres, included in larger ones, are next removed in order to reduce the total number of spheres while still covering the entire cavity space. This step is carried out repeatedly until all spheres are examined. This step significantly reduces the number of spheres while approximating the whole cavity space. If there is a water molecule present, ProCESS ignores it. FITTED will later determine whether or not the water should be considered. The grid file is outputted in mol2 format, the last column not being partial charges but the radii of the spheres.

SMART, a Tool for Ligand Preparation. We also developed the SMART module (Small Molecule Atomtyping and Rotatable Torsion assignment) which automatically identifies and labels the rotatable bonds of the ligands and assigns AMBER atom types. As the rings are not conformationally sampled in the current version of FITTED, SMART also identifies the rings²⁸ in the ligand and labels all the corresponding bonds as nonrotatable. Although no conformational sampling methods are applied to the ring, energy minimization is performed on the cyclic systems, therefore locally optimizing the ring structures. The partial charge assignment (Gasteiger–Hückel charges are recommended) is still carried out using existing software such as Sybyl.²⁹ SMART also creates reference structures of the ligands used by FITTED to compute accurate rmsds. For instance, rotamers of symmetric groups such as phenyl rings and *tert*-butyl groups are considered, creating a number of new structures that will be used as references in the atomic rmsd calculation.

FITTED 1.0, an Algorithm To Account for Protein and Ligand Flexibility. The initial proof of concept showed that the inclusion of flexibility greatly increased the accuracy of a docking run. To reduce the amount of time per run two program aspects can be investigated: removal of repetitious events (addressed by SMART and ProCESS) and increase in the quality of the individuals. The latter aspect is addressed by FITTED itself and is discussed in the following sections.

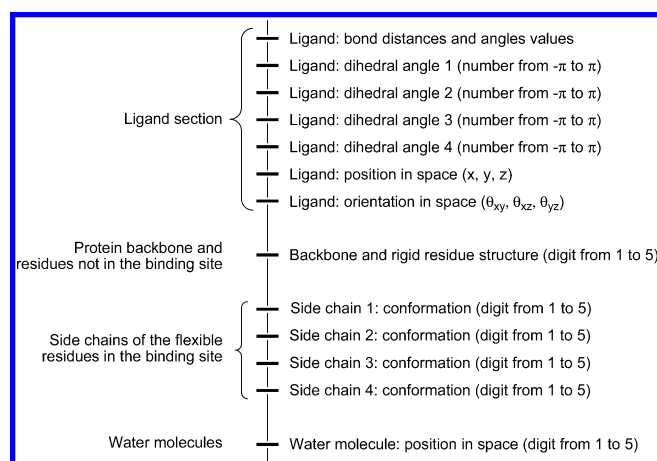


Figure 2. Chromosome describing a protein/water/ligand complex. In the illustrated case, the protein has 4 flexible residues and is represented by 5 input structures, a single water molecule is included, and the ligand has 4 rotatable bonds. A horizontal bar represents a gene.

Genetic Algorithm Implementation. Genetic algorithms (GA) have been used as optimization tools in many fields for some time. In the present work, the GA is used to optimize the binding mode. The chromosomes (Figure 2) describe the three-dimensional structure of the protein/ligand/water complex and the fitness function is the potential energy of this structure. In the illustrated case, 5 input files are used for the protein/water structures, 4 side chains are deemed flexible, and 1 bridging water molecule is considered. The first section of the chromosome codes the ligand binding mode and includes all the internal coordinates necessary to define a given conformation in a given location in space and a given orientation (often referred to as a pose). The ligand poses are therefore explicitly described in the chromosomes, and FITTED can apply a conjugate gradient algorithm to finely tune this pose.

The portion of the chromosome defining the solvated protein structure is divided into 3 sections: i. the rigid portion of the protein, including the entire backbone, ii. the side-chain conformations of the flexible binding site residues, and iii. the water molecule locations. Each side chain and rigid protein portion adopts 5 different conformations in the 5 protein input structures, referred to as call numbers with a value of 1–5. Similarly, each water molecule adopts 5 different locations again represented by a call number. Thus, libraries of side-chain conformations (1 library per side chain,

5 side-chain conformations per library), a library of 5 structures for the rigid protein portion, and a library for each water molecule (5 sets of Cartesian coordinates per water molecule) are built at the outset of the docking run from the 5 input structures. FITTED next constructs the protein/water complex from the set of digits and the libraries and adds the ligand pose to form the ternary complex. A force field energy is associated with this pose and is recalculated whenever the pose is modified.

Intelligent Design of the Initial Population. With the libraries completed, FITTED proceeds to create individuals. Each individual is first assigned a protein structure, followed by a random generation of a ligand pose.

We first thought that the required CPU time could be significantly reduced by increasing the “quality” of the initial population and focused on its generation. It is known that a population including good guesses often converges more rapidly and decreases the probability of becoming trapped in a local minimum.³⁰ In an early version, an energy threshold was used to select reasonable individuals. However, the lengthy minimization routine was used to optimize all the individuals including the many which were later discarded using this threshold, resulting in wasted CPU time. To reduce the number of unnecessary minimization steps, we envisioned the use of additional genetic operators in the form of filters (Figure 3).

A first filter was implemented that discards the poses with strong steric clashes and poses outside the protein cavity approximated by a set of spheres. If any atom of the ligand is not located within a sphere, the pose is discarded prior to potential energy evaluation. A second test is made, and only the protein/water/ligand complexes with energy below a user-defined threshold are further optimized by energy minimization.

To further improve the method, we added the possibility of exploiting experimental information by including constraints to force key interactions. For instance, ligand poses that do not interact with a given protein residue or atom (e.g., a metal) are discarded. As the grid file, the constraint file is in mol2 format. Constraints are also defined as spheres, and columns are added to the mol2 file to define the size of the constraint and its type (e.g., charge below -0.3).

Thus, in the current version, the first input protein/water file is selected, and ligand poses are randomly generated until one pose fulfills all the criteria (located within the cavity, fulfilling the constraints) (Figure 3). FITTED next constructs the complex (the corresponding chromosome) and further optimizes it through conjugate gradient energy minimization. If the optimized complex passes the last test (final energy compared to a second user-defined threshold), then it is archived. If the ligand pose does not pass, then it is discarded, and another one is generated. This procedure is reiterated with the other protein/water structures which are evenly represented in the initial population.

As expected, this implementation significantly reduced the time needed to produce a high quality initial population while including the rotational and translational degrees of freedom, which were not present in the previous method.

Evolution of Flexible Ligands. The theme of intelligent design is further carried out addressing the evolution of the ligands. The first issue addressed is the refinement of the orientation through mutation. It was observed that increasing

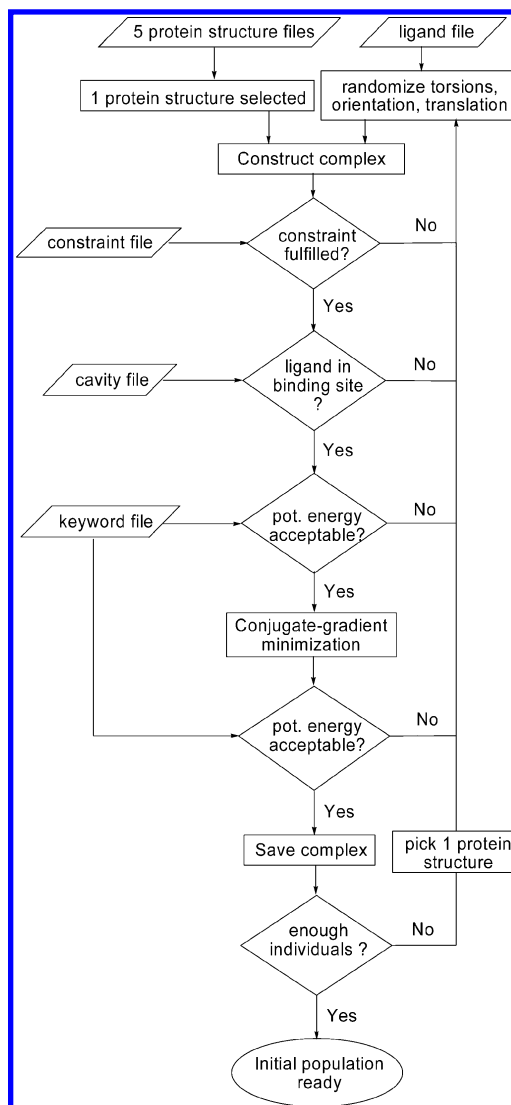


Figure 3. Generation of the initial population using a series of filters.

the probability of mutation of solely the rotational orientation of the ligand, which requires a larger sampling, resulted in an increase in the speed of convergence. Also, decreasing the range of the possible rotation mutation from 0 to 360° to $\pm 30^\circ$ yielded a better quality of individuals produced through evolution.

Second, the possibility that the best individual is further optimized without being coupled is small. To increase this possibility, we added the probability of learning. Before evolving, a small percentage of the population is further optimized by energy minimization. This approach brings the Lamarckian aspect of this GA one step further.

Evolution of Flexible Proteins – New Genetic Operators. The produced “high quality” population will then evolve using a series of genetic operators including mutations and crossover. These operators will blend the genotypes from the various individuals by swapping portions of the chromosomes (crossover) or randomly modifying genes (mutation). Parent complex structures are randomly selected from the mating pool and coupled, and children are produced by genetic operators in a steady-state way. The offspring should first pass the genotype selection described above (cavity and constraint filters) to be selected. To our knowledge this early

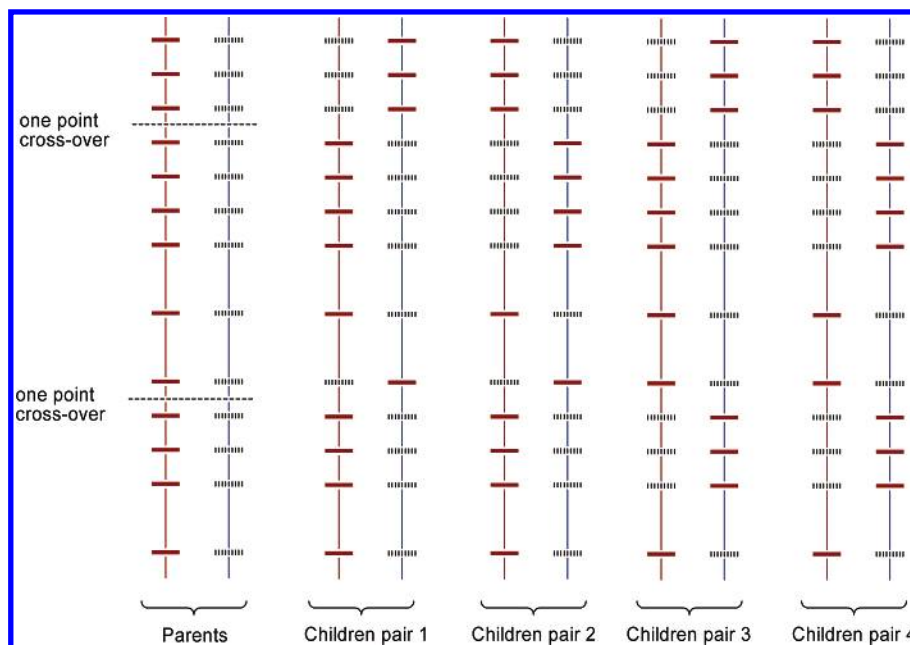


Figure 4. Four possible pairs of children generated after application of two one-point crossover operations. A horizontal bar represents a gene.

crude selection, first developed by Haupt and co-workers,³⁰ is a new concept in the GA field applied to docking methods. A proportion (user-defined) of the selection will be further optimized by energy-minimization. This energy minimization stage represents the Lamarckian aspect of the genetic algorithm. The children learn/evolve during their life (i.e., are energy-minimized) and can transmit the acquired skills to the next generation. In practice, this optimization had to be applied to a small fraction of the population. If all the structures are fully optimized, then the conformational search usually converges to high-in-energy local minima. The two best fit individuals among the parents and their children survive. This process of natural selection is based on the potential energy computed with the AMBER force field.³¹

In the current version, the input side-chain and backbone conformations and the water molecules location in space are archived in libraries, and each protein structure is described as a composite of these allowed conformations (a chromosome).²² Creating a ternary complex then requires the reconstruction of the protein structure and the addition of the water molecules and ligand. The separation between each section of the chromosome is made on purpose. FITTED will apply the genetic operators to each section independently. For instance, a single point crossover operation can be applied to the protein side-chain section, another one to the ligand internal coordinate section of the chromosome, and the last one to the water molecules (if there is more than one). As the position of the crossover is randomly selected it has a higher probability to be applied between the first and the last genes describing the ligand before the first or after the last gene. As a result, the orientation in space of the ligand (first gene of the ligand) would be somewhat linked to the backbone conformation (next gene in the chromosome). To address this artifact, when a crossover operation is performed, one of the following two options is randomly selected: the top portion of the section is kept and the bottom portion is exchanged or the top portion is exchanged and the bottom is kept. The same two options

apply to the side-chain section of the chromosome and to the water molecules (Figure 4). Crossover operations of the sections including a single gene (i.e., a single water molecule) are restricted to complete exchange or no operation. The probability to perform a crossover operation on each section is defined by the user using the appropriate keyword. Figure 4 illustrates the four possible pairs of children produced if 2 crossover operators are applied. In practice, 4 crossover operations (1 for ligand, 1 for binding site residues, 1 for the rest of the protein, and 1 for the water molecules) can be used and produce one of the 16 possibilities.

Mutation operations can also alter each gene of the chromosome except the ligand bond distances and angles. A mutation in the protein backbone, side chain, and water genes is limited to the substitution of the digit for a digit in the range defined by the number of protein input files. The mutations do not produce conformations of the protein backbone or side-chain conformations nor water molecule locations that are not in the initial libraries. As a result, FITTED will not propose protein/water structures that are not composites of the input structures. When producing a composite protein conformation, FITTED also assesses the integrity of the structure and rejects any generated protein structure that has intramolecular steric clashes.

Docking to Rigid or Flexible Proteins with FITTED.

Three options are available. First, docking to a single conformation can be performed, which allows for self- and cross-docking studies. Second, docking to a conformational ensemble can also be carried out. Using this option (referred to as "semiflexible"), the input protein structures will remain unchanged over the evolution but can be exchanged between individuals, the crossover, and mutations operating on the entire protein structures only. Third, one can use the fully flexible protein structure. With this last option, the crossover and mutation operators will be separately applied to the backbone, side chains, and water molecules.

Displaceable Water Molecules - Implementation. Bridging water molecules are often observed in protein crystal

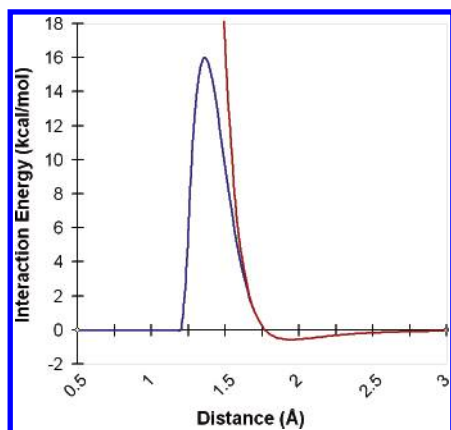


Figure 5. Interaction energy between a methanol molecule and an explicit water molecule (red) or a displaceable water molecule (blue). Cutoff distance = 1.20 Å, switching distance = 1.75 Å.

structures. To date, very few methods have been proposed to consider dynamically bound water molecules.³² We recently reported a new concept to describe displaceable water molecules.²⁰ As discussed in this previous report, the nonbonded energy function can include an energy well at an interaction distance to the water molecule but no van der Waals wall in order to simulate the water displacement. The proof-of-concept of this approach was demonstrated by using combinations of AutoDock grids modeling the “dry” and solvated RNA oligomers. The docking of aminoglycosides to these combined grids was found to be more accurate than the docking to solvated or dry RNA oligomers.²⁰ As FITTED does not make use of grids, we had to develop and implement an additional potential energy term to the AMBER function. To remove the Lennard-Jones wall for the water molecule at short distances, we introduced a switching function (SF) in the form of a scaling factor *sw* applied to the intermolecular energies involving a water molecule

$$\text{if } d < d_{\text{cutwat}}, \text{ sw} = 0.0$$

$$\text{if } d_{\text{switchwat}} < d < d_{\text{cutwat}},$$

$$\text{sw} = (d_{\text{cutwat}} - d)^2(d_{\text{cutwat}} + 2d - 3d_{\text{switchwat}})/(d_{\text{cutwat}} - d_{\text{switchwat}})^3$$

$$\text{if } d < d_{\text{switchwat}}, \text{ sw} = 1.0$$

where *sw* is the scaling factor, *d* is the shortest distance between any atom of the ligand and any atom of the water molecule, *d_{cutwat}* is the cutoff distance, and *d_{switchwat}* is the switching distance. Such functions are traditionally used to cut off long-range nonbonded interactions. In this specific case, it will be used to cut off short-range interactions.

Figure 5 represents the energy curve obtained with this new function and illustrates the interaction between methanol and a water molecule. Although the standard SFs are atom-based or group-based, this specific SF has to be molecule-based. To model a realistic situation, the water molecule should be included in the binding site (*sw* = 1 for the entire ligand) or displaced (*sw* = 0 for the entire ligand). Thus, the situations where this function ranges from 0 or 1 are artifacts. The positive energy observed in Figure 5 between 1.20 and 1.75 Å is a consequence of this function. One way

to address this issue would be to turn off the energy function as soon as it is positive. However, a continuous function between 0 and 1 was needed by the energy-minimization routine. In order to define the optimal cutoff and switching distances, the intermolecular interaction within complexes such as methanol–water and *N*-methyl acetamide–water was investigated. In all cases, the interaction energy between the molecules was positive at distances below 1.75 Å, selected as the optimal switching distance. Therefore, this SF applies only when the interaction energy with the water molecule is repulsive. The SF reached a maximum of about 15 kcal/mol when a cutoff distance of 1.20 Å was used (see Figure 5).

Applied to the docking of molecules, this potential energy function penalizes the poses that do not interact favorably with a water molecule (distance < cutoff distance = 1.75 Å) nor displace it completely (distance > switching distance = 1.20 Å) and will consequently favor ligand poses that either interact or fully displace the water molecules.

Displaceable Water Molecules - Optimizing Water Evolution. In the present work, critical water molecules are either maintained when present in the crystal structures or added by analogy to other structures and their orientation optimized by energy minimization when missing. Initial attempts have shown that the prediction of the occurrence of water molecules in the complexes was not accurate. In practice, we observed that the ligand pose was first optimized (with greater decrease of the total potential energy), followed by the refinement of the protein structure and finally the water molecules. However, most of the water location possibilities (one per protein structure) have been removed throughout the generations. We then found that higher mutation rates increased the accuracy by increasing the sampling of the water molecules. In order to address this issue, we implemented a ramping mutation rate for the water molecules. This ramping is achieved by using a quadratic function.

$$p_{\text{mutwat}} = p_{\text{max}} (nth \text{ gen.}/\text{max. number of gens.})^4$$

Thus, very low mutation rates are applied at the early stages of the evolution, while larger rates are used at the late stages. One drawback of the use of the AMBER force field is the lack of directionality of the hydrogen bond term. Evaluation of the free energy of binding of the water molecules was also a concern. In the current version of FITTED, water molecules are considered as part of the protein. To account at least partly for the entropy cost associated with the capture of a water molecule, a penalty is added to the final score whenever a water molecule is maintained. This number is arbitrary as this penalty is system-dependent and should also include the enthalpic contribution to the binding of the water. Work is in progress in our laboratory to include directional hydrogen bonds in the next version of FITTED and to improve the scoring of the free energy of binding of the water molecules to protein/ligand complexes.

Scoring Function. The AMBER force field was implemented in FITTED and used during the actual docking with a higher weight for the intermolecular interactions than for the internal energy. Very few scoring functions include a

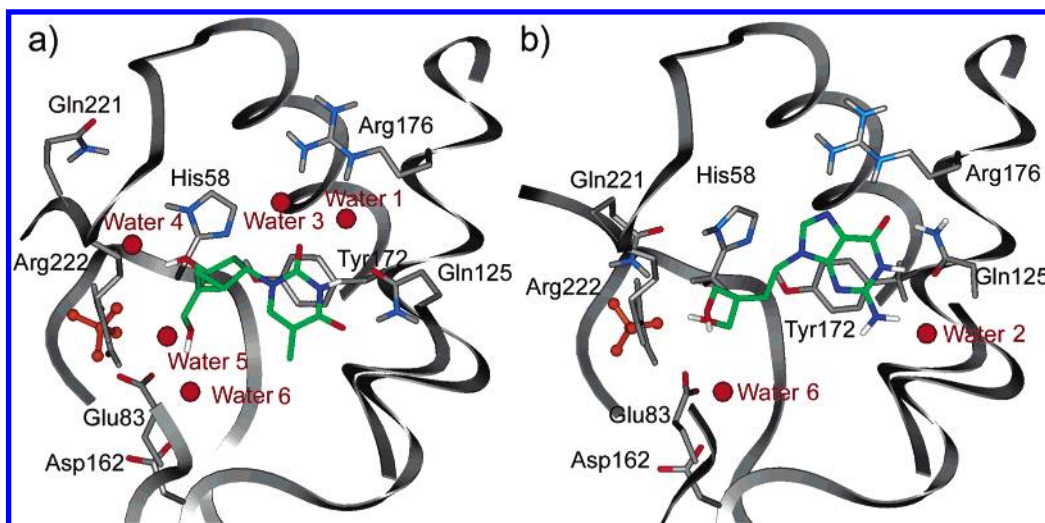


Figure 6. Bridging water molecules and flexible binding site residues in TK/inhibitor complexes: (a) 1e2k and (b) 1ki3. Cocrystallized sulfate is shown in orange.

term accounting for the protein entropy. In the present case, using a force field does not permit the evaluation of the entropy contribution to the free energy of binding. Understanding that the mobility and entropy of flexible residues are modulated by the ligands, we have proposed to estimate the free energy of binding to the flexible residues.²² First, the stronger the interactions are, the tighter a ligand is bound. Then, the tighter a ligand is bound, the more frozen the surrounding side chains are. We proposed to account for this entropy/enthalpy compensation by reducing the interactions with flexible residues. In practice, the interaction with flexible side chains was scaled down by the use of a new set of atom types and partial charges²² assigned by ProCESS. The final poses were then scored using our scoring function Rank-Score²² also implemented in FITTED.

RESULTS AND DISCUSSION

Selection of the Testing Set. As FITTED incorporates protein flexibility and displaceable water molecules, a selection of protein/inhibitor complexes should be made to evaluate these aspects of FITTED. The selected complexes are listed as Supporting Information.

As HIV-1 protease (HIVP)/inhibitor complexes often involve a bridging water molecule, this enzyme was selected as a first test case. Although HIVP is a flexible protein, the inhibitors usually bind to the closed form, and HIVP is not considered a highly flexible protein in docking studies. However, slight adjustments were observed, and rmsds of 0.5–1.4 Å between binding site residues were computed. HIVP can exhibit two different protonation states, either one or both catalytic aspartic acid side chains being protonated.³³ In most cases, inhibitors binding to the catalytic dyad via a diol moiety favor the diprotonation of the catalytic aspartates, while monoalcohols or other functional groups favor the monoprotonated state. We therefore decided to prepare two sets of protein files: 1b6l, 1eby, 1hpo, 1hvp, and 1pro protein structures were monoprotonated as discussed in the experimental section, while 1ajv, 1ajx, 1hvr, 1hwr, and 1qbs were diprotonated.³³ Only the crystal structures 1b6l, 1eby, and 1hvp featured a bridging water molecule. This same water molecule was therefore added to the other 7 protein structures

to allow FITTED to select whether or not this water is needed.

Similarly, thymidine kinase (TK) is a flexible protein, and inhibitors often bind experiencing interactions with the protein relayed by many water molecules. Interestingly, a first water (water molecules 1 and 2 in Figure 6) can be located at two different positions following Gln125 side-chain conformational changes. The combined water displacement/Gln side-chain flip will be investigated in great detail. As illustrated in Figure 6, either the Gln125 carbonyl oxygen (Figure 6a) or the amide hydrogens (Figure 6b) point toward the Arg176 side chain. The first Gln125 side-chain conformation shown in Figure 6a is observed in 1e2k, 1e2p, 1ki4, 1ki8, and 1of1, while the second conformation is observed in 1ki3, 1ki7, 2ki5, and 1qhi. Similarly, Water 4 can be displaced by Gln221. These two enzymes (HIVP and TK) together with oligopeptide binding protein A (OppA) were also selected as test cases by the GOLD developers in order to evaluate the reliability of their method accounting for bridging water molecules.³² However, Verdonk et al. considered three water molecules interacting with the nucleoside base of the TK inhibitors, whereas we considered another three interacting with the ribose part of these inhibitors, for a total of six water molecules. As illustrated in Figure 3, these six waters participate in multiple hydrogen bonds with both the ligands and the proteins. As for HIVP, missing waters were added by analogy with crystal structures featuring these waters.

Factor Xa (FXa) and its homolog trypsin were also included in the validation set. FXa/inhibitor complexes show from 0–2 water molecules involved in the ligand binding, while a single bridging water molecule is observed in the selected trypsin/inhibitor complexes (Figure 7). The first water molecule interacts with both the inhibitor cationic moieties and the protein backbone (Ile227 in FXa and Val205 in trypsin), while the second one bridges the inhibitor cation with the key Asp189 side chain of factor Xa. The specific shape of these two deep binding sites featuring a narrow pocket made the conformational sampling problematic. Thus, a larger population size (200) was used in order to reach convergence.

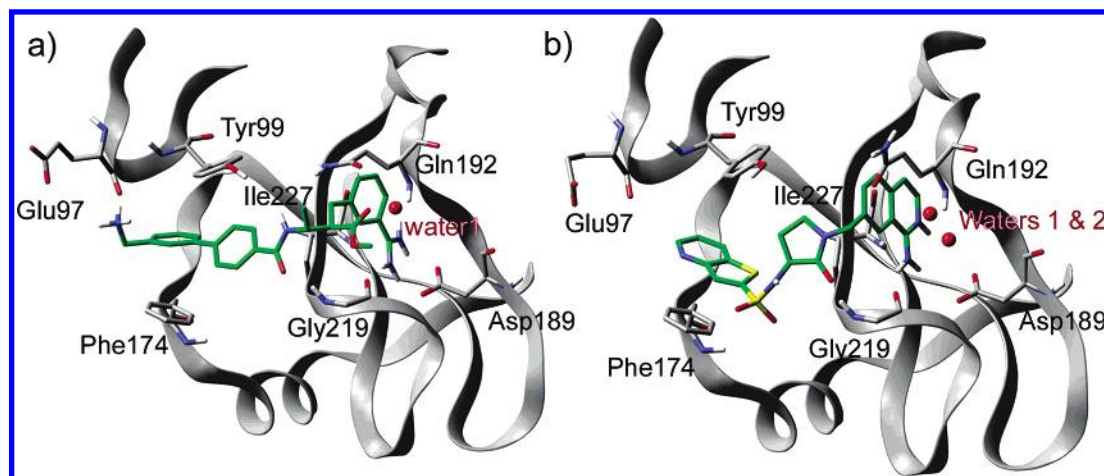


Figure 7. Bridging water molecules and flexible binding site residues in FXa/inhibitor complexes: (a) 1ezq and (b) 1f0r.

Table 1. Self-Docking – HIV-1 Protease Inhibitors

| protonation Asp25 | PDB code | obsd water ^c | docking to protein ^a | | docking to protein + water molecule ^b | | |
|-------------------|----------|-------------------------|---------------------------------|--------------------|--|-------------------------|--------------------|
| | | | ligand ^d | score ^e | Ligand ^d | pred water ^f | score ^e |
| monoprotonated | 1b6l | 1 | 1.10 | −9.8 | 0.55 | 1 | −11.6 |
| | 1eby | 1 | 2.68 | −9.3 | 4.55 | 0 | −9.4 |
| | 1hpo | 0 | 2.29 | −9.1 | 0.94 | 0 | −10.1 |
| | 1hpu | 1 | 1.02 | −8.7 | 1.19 | 1 | −8.5 |
| | 1pro | 0 | 0.82 | −5.2 | 0.72 | 0 | −5.7 |
| diprotonated | 1ajv | 0 | 0.91 ^g | −10.0 | 0.59 | 0 | −11.4 |
| | 1ajx | 0 | 0.82 | −9.4 | 0.77 | 0 | −9.9 |
| | 1hvr | 0 | 0.49 | −11.7 | 0.40 | 0 | −12.2 |
| | 1hwr | 0 | 0.60 | −7.5 | 0.52 | 0 | −8.1 |
| | 1qbs | 0 | 5.05 | −7.5 | 5.14 | 0 | −8.2 |

^a Water molecules removed prior to docking. ^b Water molecule known as Water 301 was retained, and the function describing the interaction between ligand and water molecules is applied. ^c Water molecule as proposed by FITTED; 1 and 0 define the presence or absence of the water molecule, respectively. ^d rmsd (in Å): criterion of success of 2.0 Å. ^e Score in arbitrary units. ^f Water molecules as proposed by FITTED. Bold numbers highlight failures.

Table 2. Self-Docking – Thymidine Kinase Inhibitors

| | | | | | | | docking to protein ^a | | docking to protein + water molecule ^b | | | | | | |
|-------------|---|---|---|---|---|---|---------------------------------|--------------------|--|-----------------------------------|---|----------|----------|----------|--------------------|
| | | | | | | | ligand ^d | score ^e | ligand ^d | pred water molecules ^f | | | | | score ^e |
| 1e2k | 1 | 0 | 1 | 1 | 1 | 1 | 0.63 | −6.1 | 0.66 | 1 | 0 | 1 | 1 | 0 | −7.1 |
| 1e2p | 1 | 0 | 1 | 1 | 1 | 1 | 2.69 ^g | −4.7 | 2.03 ^g | 1 | 0 | 1 | 1 | 1 | −5.2 |
| 1ki3 | 0 | 1 | 0 | 0 | 0 | 1 | 1.86 | −5.9 | 1.84 | 0 | 1 | 0 | 1 | 1 | −6.1 |
| 1ki4 | 1 | 0 | 1 | 1 | 1 | 1 | 0.43 | −6.9 | 0.66 | 1 | 0 | 1 | 1 | 1 | −7.7 |
| 1ki7 | 1 | 0 | 1 | 1 | 1 | 0 | 5.79 | −5.1 | 5.76 | 1 | 0 | 1 | 1 | 1 | −4.8 |
| 1ki8 | 1 | 0 | 1 | 1 | 1 | 0 | 0.77 | −6.2 | 0.64 | 1 | 0 | 1 | 1 | 1 | −6.9 |
| 2ki5 | 0 | 1 | 1 | 1 | 1 | 1 | 1.10 | −5.5 | 0.45 | 0 | 1 | 0 | 1 | 1 | −6.3 |
| 1of1 | 1 | 0 | 1 | 1 | 1 | 1 | 0.37 | −6.1 | 0.29 | 1 | 0 | 1 | 1 | 1 | −6.8 |
| 1qhi | 0 | 1 | 1 | 1 | 1 | 0 | 0.47 | −7.2 | 0.66 | 0 | 1 | 0 | 1 | 1 | −7.8 |

^a Water molecules removed prior to docking. ^b Two to 6 water molecules (see text) were retained, and the function describing the interaction between ligand and water molecules is applied. ^c Water molecules observed or not in crystal structures: 1 and 0 define the presence or absence of each water molecule, respectively. ^d rmsd (in Å): criterion of success of 2.0 Å. ^e Score in arbitrary units. ^f Water molecules as proposed by FITTED. ^g When considering the prochiral nature of 1e2p ligand, these rmsds were equivalent to rmsds below 1.0 Å (see text). Bold numbers highlight failures.

We completed this validation study with a small set of metalloprotease (MMP-3, stromelysin-1) inhibitors for a total of 33 complexes. Most of the known MMP inhibitors chelate the catalytic zinc cation and are of interest to evaluate the accuracy of FITTED to reproduce the metal ligation. As the metal chelation is a short-range interaction, we implemented a specific term using a potential similar to the LJ12-10 used for hydrogen bonds.

All these protein structures were processed using ProCESS (a typical keyword file is given as Supporting Information) prior to their use with FITTED, and the ligands were prepared with SMART.

Docking Using FITTED 1.0 – General Considerations.

All the compounds were first self-docked to their corresponding protein structure in the presence or absence of water. Tables 1–3 summarize the data obtained for these docking studies. This first set of docking runs was carried out to evaluate the impact of the new potential energy term for the displaceable water molecules. A cross docking study was next carried out to evaluate the impact of the protein structure on the docking accuracy (Tables 4–6). These same ligands were next docked to the “semiflexible” proteins and to the fully flexible proteins in order to evaluate the ability of FITTED to predict the protein structure (Tables 7–9).

Table 3. Self-Docking – Factor Xa Trypsin and MMP-3 Inhibitors

| protein | PDB code | obsd water ^c | | docking to protein ^a | | docking to protein + water molecule ^b | | | |
|-----------|----------|-------------------------|---|---------------------------------|--------------------|--|-------------------------|----------|--------------------|
| | | | | ligand ^d | score ^e | ligand ^d | pred water ^f | | score ^e |
| factor Xa | 1ezq | 1 | 0 | 3.32 | −7.4 | 0.82 | 0 | 0 | −11.1 |
| | 1f0r | 1 | 1 | 2.33 | −8.3 | 1.88 | 0 | 0 | −8.1 |
| | 1fjs | 1 | 0 | 3.64 | −7.7 | 1.78 | 0 | 0 | −8.8 |
| | 1nfu | 0 | 0 | 2.57 | −8.9 | 1.50 | 0 | 0 | −8.0 |
| | 1xka | 1 | 0 | 1.13 | −8.3 | 0.87 | 0 | 0 | −8.4 |
| trypsin | 1f0u | 1 | - | 2.92 | −5.9 | 3.95 | 1 | - | −6.7 |
| | 1o2j | 1 | - | 1.03 | −5.9 | 0.94 | 1 | - | −6.2 |
| | 1o3g | 1 | - | 1.35 | −6.9 | 1.69 | 1 | - | −7.3 |
| | 1o3i | 1 | - | 0.70 | −6.5 | 0.68 | 1 | - | −6.7 |
| | 1qbo | 1 | - | 3.84 | −7.6 | 3.49 | 1 | - | −6.8 |
| MMP-3 | 1b8y | - | - | 1.15 | −9.4 | - | - | - | - |
| | 1bwi | - | - | 6.35 | −5.8 | - | - | - | - |
| | 1ciz | - | - | 1.22 | −10.6 | - | - | - | - |
| | 1d8m | - | - | 2.99 | −6.0 | - | - | - | - |

^a Water molecules removed prior to docking. ^b Zero to 2 water molecules (see text) were retained, and the function describing the interaction between ligand and water molecules is applied. ^c Water molecules observed or not in crystal structures: 1 and 0 define the presence or absence of each water molecule, respectively. ^d rmsd (in Å): criterion of success of 2.0 Å. ^e Score in arbitrary units. ^f Water molecules as proposed by FITTED. Bold numbers highlight failures.

Table 4. Cross-Docking and Docking to Multiple Conformations – HIV-1 Protease Inhibitors

| | docking to rigid proteins | | | | | statistics for the best scoring pose ^a | | | |
|------|---------------------------|-------------|-------------|-------------|-------------|---|----------------------|--------------------|--------------------|
| | 1b6l | 1eby | 1hpo | 1hvp | 1pro | ligand ^b | protein ^c | water ^d | score ^e |
| 1b6l | 0.55 | 0.83 | 3.37 | 1.11 | 1.04 | 0.55 | 0.00 | 1 | −11.6 |
| 1eby | 2.57 | 4.55 | 2.86 | 6.15 | 2.72 | 2.86 | 0.96 | 1 | −8.9 |
| 1hpo | 4.64 | 3.62 | 0.94 | 4.26 | 2.4 | 0.94 | 0.00 | 0 | −10.1 |
| 1hvp | 4.09 | 3.39 | 2.01 | 1.19 | 3.53 | 2.01 | 1.00 | 0 | −8.8 |
| 1pro | 0.62 | 1.01 | 0.86 | 0.78 | 0.72 | 0.72 | 0.00 | 0 | −5.7 |

| | docking to rigid proteins | | | | | statistics for the best scoring pose ^a | | | |
|------|---------------------------|------|------|-------------|-------------|---|----------------------|--------------------|--------------------|
| | 1ajv | 1ajx | 1hvr | 1hwr | 1qbs | ligand ^b | protein ^c | water ^d | score ^e |
| 1ajv | 0.59 | 1.26 | 1.46 | 1.52 | 1.12 | 1.12 | 0.87 | 0 | −10.1 |
| 1ajx | 0.73 | 0.77 | 1.1 | 0.75 | 0.73 | 0.73 | 0.81 | 0 | −9.6 |
| 1hvr | 1.9 | 1.27 | 0.4 | 1.22 | 0.77 | 0.77 | 0.72 | 0 | −11.9 |
| 1hwr | 0.68 | 0.78 | 0.85 | 0.52 | 0.67 | 0.78 | 0.84 | 0 | −8.9 |
| 1qbs | 5.35 | 1.49 | 1.17 | 5.11 | 5.15 | 1.17 | 0.72 | 0 | −10.3 |

^a Each ligand was docked to the 5 protein structure, and the best scoring of the 5 final poses was selected. ^b rmsd (in Å): criterion of success of 2.0 Å. ^c rmsd (in Å): criterion of success: better than average rmsd; average rmsd between protein structures computed on the binding site residues: 0.91 Å for the first five structures (one Asp 25 protonated) and 0.77 Å for the last five structures (AspA25 and AspB25 protonated). ^d Water molecules as proposed by FITTED; 1 and 0 define the presence or absence of each water molecule, respectively. ^e Score in arbitrary units. Bold numbers highlight failures.

For the statistical analysis, we considered that the ligand pose was accurately predicted when the atomic rmsd relative to the crystal structure was below 2.0 Å and that the protein structure prediction was correct when the rmsd was below the average rmsd between the series of protein structures used as input (i.e., when the prediction is better than a protein structure randomly chosen). Finally, we considered the water molecules to be accurately predicted when the occurrence was right. A set of 10 runs was carried out for each inhibitor, in order to demonstrate the convergence of the protocol. In most cases at least 5 out of the 10 runs led to similar poses (difference between computed rmsds below 0.5 Å). Although 100 individuals were enough for the docking of thymidine kinase inhibitors, a larger initial population (200) was required for the other 4 proteins in order to reach the convergence criterion.

Self-Docking and Effect of Displaceable Waters. Among the 5 proteins investigated, 4 proteins can bind ligands through one or more bridging water molecules. In order to evaluate the impact on the docking accuracy of the potential energy term developed for the water molecules, two sets of experiments were carried out. In the first set of experiments, the water molecules were removed from the protein structures, and inhibitors were docked back to their corresponding protein structure (self-docking). In the second set of experiments, the water molecules were kept, and the developed potential energy term for the water molecule was used.

Table 1 presents the results of the self-docking study for HIVP inhibitors. As can be seen in the third and fifth columns, 7 (without water) or 8 (with water) out of the 10 inhibitors were self-docked within 1.2 Å from the experimentally observed binding modes. Interestingly, 1hpo was properly docked when the displaceable water potential was used and the water is predicted to be displaced. However, when the water was removed prior to docking, 1hpo, which is known to displace the water molecule, is misdocked. We attribute this unexpected result to the energy hill shown in Figure 5. As postulated above, this energy potential tends to favor either the complete displacement of the water or favorable interactions with the water while disfavoring intermediate docked poses as the one proposed when the water is removed prior to docking. A close look at Table 1 also revealed that the rmsds are systematically higher when the water is removed prior to docking. We believe that this can be attributed to the energy potential used to model the displaceable water molecules.

The TK inhibitors were next docked to the rigid protein in self-docking experiments (Table 2). Seven out of nine inhibitors were docked with rmsd below 2.0 Å, with 6 inhibitors self-docked within 1.1 Å of the observed binding mode. A special situation arose with prochiral compound 1e2p. As shown in Figure 8, the rmsd of 2.03 Å computed for the docked pose when the water molecules were considered was attributed to the exchange of C-1 and C-2 groups. Considering the two methylenol groups as equivalent reduces the rmsd to 0.77 Å. Docking of 1e2p was therefore

Table 5. Cross-Docking and Docking to Multiple Conformations – Thymidine Kinase Inhibitors

| | docking to rigid proteins | | | | | | | | | statistics for the best scoring pose ^a | | | | | | | | |
|------|---------------------------|-------------------------|-------------|-------------|-------------|-------------|-------------------------|-------------|-------------------------|---|----------------------|--------------------|----------|----------|----------|--------------------|----------|------|
| | 1e2k | 1e2p | 1ki3 | 1ki4 | ki7 | 1ki8 | 2ki5 | 1of1 | 1qhi | ligand ^b | protein ^c | water ^d | | | | score ^e | | |
| 1e2k | 0.66 | 2.11 | 3.42 | 0.76 | 0.83 | 0.84 | 0.96 | 0.78 | 1.31 | 0.66 | 0.00 | 1 | 0 | 1 | 1 | 0 | 1 | −8.1 |
| 1e2p | 2.24^f | 2.03^f | 0.97 | 0.74 | 0.78 | 1.41 | 2.75^f | 1.20 | 2.03^f | 2.24^f | 0.59 | 1 | 0 | 0 | 1 | 1 | 1 | −6.2 |
| 1ki3 | 2.36 | 2.62 | 1.84 | 2.43 | 2.25 | 2.50 | 2.61 | 2.94 | 1.74 | 1.74 | 0.78 | 1 | 0 | 1 | 1 | 1 | 1 | −6.1 |
| 1ki4 | 2.48 | 2.36 | 3.38 | 0.66 | 1.11 | 0.73 | 2.35 | 2.52 | 1.00 | 0.66 | 0.00 | 1 | 0 | 1 | 1 | 1 | 1 | −8.2 |
| 1ki7 | 5.79 | 5.67 | 5.55 | 5.08 | 5.76 | 5.25 | 5.15 | 5.65 | 5.67 | 5.67 | 0.87 | 0 | 0 | 0 | 1 | 1 | 1 | −5.8 |
| 1ki8 | 3.8 | 3.93 | 2.35 | 1.91 | 1.19 | 0.64 | 3.92 | 3.84 | 1.29 | 0.64 | 0.00 | 1 | 0 | 1 | 1 | 1 | 1 | −7.4 |
| 2ki5 | 2.29 | 3.22 | 1.28 | 2.13 | 2.08 | 1.90 | 0.45 | 2.22 | 1.21 | 1.90 | 1.11 | 1 | 0 | 1 | 1 | 1 | 1 | −6.7 |
| 1of1 | 0.39 | 0.49 | 1.14 | 0.60 | 0.78 | 0.89 | 0.49 | 0.29 | 0.81 | 0.29 | 0.00 | 1 | 0 | 0 | 1 | 1 | 1 | −7.3 |
| 1qhi | 2.41 | 2.29 | 1.18 | 5.43 | 1.68 | 2.13 | 1.10 | 2.19 | 0.67 | 0.67 | 0.00 | 0 | 1 | 1 | 1 | 1 | 1 | −6.1 |

^a Each ligand was docked to the 5 protein structure, and the best scoring of the 5 final poses was selected. ^b rmsd (in Å): criterion of success of 2.0 Å. ^c rmsd (in Å): criterion of success: better than average rmsd; average rmsd between protein structures computed on the binding site residues: 0.92 Å. ^d Water molecules as proposed by FITTED; 1 and 0 define the presence or absence of each water molecule, respectively. ^e Score in arbitrary units. ^f Equivalent to rmsds below 1.0 Å if the prochiral nature of the ligand is considered. Bold numbers highlight failures.

Table 6. Cross-Docking and Docking to Multiple Conformations – Factor Xa, Trypsin and MMP-3 Inhibitors

| | docking to rigid proteins | | | | | statistics for the best scoring pose ^a | | | |
|------|---------------------------|-------------|-------------|-------------|-------------|---|----------------------|--------------------|--------------------|
| | 1ezq | 1f0r | 1fsj | 1nfu | 1xka | ligand ^b | protein ^c | water ^d | score ^e |
| 1ezq | 0.82 | 9.14 | 3.53 | 3.45 | 4.8 | 0.82 | 0.00 | 0 | −11.1 |
| 1f0r | 2.42 | 1.89 | 2.31 | 2.49 | 2.42 | 2.49 | 0.75 | 0 | −9.4 |
| 1fsj | 3.3 | 2.81 | 1.78 | 2.24 | 3.22 | 1.78 | 0.00 | 0 | −8.8 |
| 1nfu | 2.05 | 2.17 | 2.26 | 1.5 | 3.79 | 1.5 | 0.00 | 0 | −9.7 |
| 1xka | 1.66 | 1.5 | 1.64 | 1.58 | 0.87 | 0.87 | 0.00 | 0 | −8.4 |

| | docking to rigid proteins | | | | | statistics for the best scoring pose ^a | | | |
|------|---------------------------|-------------|-------------|-------------|-------------|---|----------------------|--------------------|--------------------|
| | 1f0u | 1o2j | 1o3g | 1o3i | 1qbo | ligand ^b | protein ^c | water ^d | score ^e |
| 1f0u | 3.95 | 4.21 | 2.16 | 4.98 | 5.50 | 5.50 | 0.74 | 1 | −6.9 |
| 1o2j | 1.33 | 0.94 | 0.80 | 4.16 | 1.43 | 0.80 | 0.55 | 1 | −6.3 |
| 1o3g | 0.59 | 0.79 | 1.69 | 0.67 | 1.22 | 0.67 | 0.31 | 1 | −7.6 |
| 1o3i | 0.59 | 0.93 | 0.94 | 0.69 | 1.06 | 0.69 | 0.00 | 1 | −6.7 |
| 1qbo | 5.23 | 4.14 | 3.89 | 4.30 | 3.49 | 3.89 | 1.09 | 1 | −7.4 |

| | docking to rigid proteins | | | | statistics for the best scoring pose ^a | | |
|------|---------------------------|-------------|-------------|-------------|---|----------------------|--------------------|
| | 1b8y | 1bwi | 1ciz | 1d8m | ligand ^b | protein ^c | score ^e |
| 1b8y | 1.15 | 1.51 | 1.38 | 2.30 | 1.15 | 0.00 | −9.4 |
| 1bwi | 5.64 | 6.35 | 8.95 | 6.40 | 6.35 | 0.00 | −5.8 |
| 1ciz | 1.15 | 4.53 | 1.23 | 4.33 | 1.15 | 0.45 | −10.1 |
| 1d8m | 1.22 | 6.23 | 2.21 | 2.99 | 1.22 | 1.11 | −7.3 |

^a Each ligand was docked to the 5 protein structure, and the best scoring of the 5 final poses was selected. ^b rmsd (in Å): criterion of success of 2.0 Å. ^c rmsd (in Å): criterion of success: better than average rmsd; average rmsd between protein structures computed on the binding site residues: factor Xa: 0.86 Å, trypsin: 0.90 Å, MMP-3: 0.92 Å. ^d Water molecules as proposed by FITTED; 1 and 0 define the presence or absence of each water molecule, respectively. ^e Score in arbitrary units. Bold numbers highlight failures.

considered as successful, raising the success rate to 8 out of 9 complexes.

Surprisingly, even though many water molecules are involved in the ligand/protein complexes, the docking to the “dry protein” was as accurate as the docking to the solvated protein. Only the docking of 2ki5 was slightly affected by the absence of water. The ten runs carried out with 1ki7 were consistently leading to the same wrong conformation. In this case, the ribose ring and the base of the nucleotide mimics were inverted within the binding site.

We next investigated the two sets of charged trypsin and FXa inhibitors. In this case, the need for water molecules was clear (Table 3). Without any water molecules FITTED docked only 4 out of the 10 inhibitors properly, while 8

Table 7. Docking to Flexible Proteins - HIV-1 Protease Inhibitors

| | docking to semiflexible protein | | | | docking to fully flexible protein | | | |
|------|---------------------------------|----------------------|--------------------|--------------------|-----------------------------------|----------------------|----------|-------|
| | ligand ^a | protein ^b | water ^c | score ^d | ligand ^a | protein ^b | water | score |
| 1b6l | 1.06 | 0.00 | 1 | −11.0 | 1.08 | 0.53 | 1 | −11.4 |
| 1eby | 3.62 | 0.85 | 0 | −9.3 | 6.06 | 1.02 | 0 | −8.6 |
| 1hpo | 4.03 | 0.99 | 0 | −8.1 | 3.25 | 1.16 | 0 | −8.5 |
| 1hvp | 3.88 | 1.00 | 1 | −10.3 | 1.54 | 0.79 | 1 | −10.0 |
| 1pro | 0.51 | 0.00 | 0 | −5.9 | 0.94 | 0.59 | 0 | −5.6 |
| 1ajv | 0.75 | 0.00 | 0 | −11.4 | 1.46 | 1.02 | 0 | −10.6 |
| 1ajx | 0.85 | 0.84 | 0 | −9.3 | 1.77 | 0.72 | 0 | −10.0 |
| 1hvr | 1.72 | 0.00 | 0 | −9.7 | 1.59 | 0.67 | 0 | −11.6 |
| 1hwr | 0.79 | 0.81 | 0 | −8.8 | 0.58 | 0.71 | 0 | −8.9 |
| 1qbs | 1.22 | 0.72 | 0 | −11.0 | 1.32 | 0.59 | 0 | −11.0 |

^a rmsd (in Å): criterion of success of 2.0 Å. ^b rmsd (in Å): criterion of success: better than average rmsd; average rmsd between protein structures computed on the binding site residues: 0.91 Å for the first five structures (one Asp 25 protonated) and 0.77 Å for the last five structures (AspA25 and AspB25 protonated). ^c Water molecules as proposed by FITTED; 1 and 0 define the presence or absence of the water molecule, respectively. ^d Score in arbitrary units. Bold numbers highlight failures.

inhibitors were accurately docked when the water molecules were considered. A close look at the failures does not reveal any major mistake. For instance, the proposed poses for 1f0u were interacting with Asp171 as experimentally observed. However, the hydrophobic biaryl moiety of 1f0u was not located in the correct pocket, with the terminal ammonium group forming a hydrogen bond with Tyr76 instead of Asn79. The surprise comes from the water prediction. In three cases (1ezq, 1f0r, and 1fsj) the occurrence of water molecules is not accurately predicted, but nevertheless the inhibitors are docked better than without the water molecules.

In contrast, the occurrence of water is accurately predicted when docking trypsin inhibitors and the removal of the water do not affect the docking.

MMP inhibitors were docked with low accuracy, with only 2 inhibitors having rmsd < 2.0 Å. This small set is clearly not large enough to fully assess FITTED for metalloenzymes.

Overall, these first experiments demonstrated the ability of FITTED to fully sample the ligand conformational space and assign better scores to experimentally observed poses. This first study also validated the water molecule prediction method since the occurrence of the so-called water 301 in HIVP and water molecules in TK and trypsin is right in most cases. Unexpectedly, this additional energy term also helps in the docking of inhibitors that displace water.

Table 8. Docking to Flexible Proteins - Thymidine Kinase Inhibitors

| | docking to semiflexible protein | | | | | | | | | docking to fully flexible protein | | | | | | | | |
|------|---------------------------------|----------------------|---------------------------------------|----------|----------|---|----------|----------|--------------------|-----------------------------------|----------------------|---------------------------------------|----------|----------|----------|----------|----------|--------------------|
| | ligand ^a | protein ^b | occurrence of water mol. ^c | | | | | | score ^d | ligand ^a | protein ^b | occurrence of water mol. ^c | | | | | | score ^d |
| 1e2k | 0.67 | 0.00 | 1 | 0 | 1 | 1 | 0 | 1 | −7.0 | 0.75 | 0.61 | 1 | 0 | 1 | 1 | 0 | 1 | −7.2 |
| 1e2p | 0.51 | 0.88 | 1 | 0 | 0 | 1 | 0 | 1 | −5.8 | 0.95 | 0.93 | 1 | 0 | 1 | 1 | 1 | 1 | −5.7 |
| 1ki3 | 1.46 | 0.00 | 0 | 1 | 1 | 0 | 0 | 1 | −6.8 | 1.35 | 0.90 | 0 | 0 | 0 | 1 | 1 | 1 | −6.8 |
| 1ki4 | 0.64 | 0.00 | 1 | 0 | 1 | 1 | 1 | 1 | −7.3 | 0.77 | 0.89 | 1 | 0 | 0 | 1 | 1 | 1 | −7.9 |
| 1ki7 | 5.20 | 1.01 | 1 | 0 | 1 | 1 | 1 | 1 | −5.4 | 5.25 | 1.11 | 0 | 1 | 0 | 1 | 1 | 0 | −6.5 |
| 1ki8 | 0.60 | 0.96 | 1 | 0 | 1 | 1 | 1 | 1 | −6.7 | 0.65 | 0.53 | 1 | 0 | 1 | 1 | 1 | 1 | −7.8 |
| 2ki5 | 1.92 | 0.89 | 0 | 1 | 1 | 1 | 1 | 1 | −5.6 | 1.62 | 0.80 | 1 | 0 | 1 | 1 | 1 | 1 | −6.8 |
| 1of1 | 0.35 | 0.26 | 1 | 0 | 1 | 1 | 1 | 1 | −6.7 | 0.75 | 0.99 | 1 | 0 | 1 | 1 | 1 | 1 | −7.3 |
| 1qhi | 0.96 | 0.00 | 1 | 0 | 0 | 1 | 1 | 0 | −7.5 | 0.64 | 0.69 | 0 | 1 | 0 | 1 | 1 | 1 | −8.2 |

^a rmsd (in Å): criterion of success of 2.0 Å. ^b rmsd (in Å): criterion of success: better than average rmsd; average rmsd between protein structures computed on the binding site residues: 0.92 Å. ^c Water molecules as proposed by FITTED; 1 and 0 define the presence or absence of each water molecule, respectively. ^d Score in arbitrary units. Bold numbers highlight failures.

Table 9. Cross-Docking and Docking to Flexible Proteins – Factor Xa, Trypsin, and MMP-3 Inhibitors

| | docking to semiflexible protein | | | | docking to fully flexible protein | | | | | |
|------|---------------------------------|----------------------|--------------------|--------------------|-----------------------------------|-------------------------|--------------------|--------------------|----------|-------|
| | ligand ^a | protein ^b | water ^c | score ^d | ligand ^a | protein ^b | water ^c | score ^d | | |
| 1ezq | 1.34 | 0.00 | 1 | 0 | -10.2 | 9.64^e | 0.92 | 1 | 0 | -8.5 |
| 1f0r | 2.50^d | 0.75 | 0 | 0 | -8.1 | 2.32^f | 0.63 | 0 | 0 | -9.7 |
| 1fjs | 2.45 | 0.77 | 0 | 0 | -9.1 | 3.24 | 1.11 | 1 | 0 | -8.6 |
| 1nfu | 1.87 | 0.70 | 0 | 0 | -8.7 | 1.17 | 0.71 | 0 | 0 | -9.6 |
| 1xka | 1.31 | 0.91 | 0 | 0 | -8.2 | 1.52 | 0.70 | 0 | 0 | -8.7 |
| 1f0u | 6.11 | 1.04 | 1 | - | -6.7 | 4.25 | 0.87 | 1 | - | -7.6 |
| 1o2j | 1.06 | 0.33 | 1 | - | -6.5 | 1.30 | 0.58 | 1 | - | -7.1 |
| 1o3g | 0.83 | 0.66 | 1 | - | -7.2 | 0.82 | 0.77 | 1 | - | -7.7 |
| 1o3i | 1.24 | 1.32 | 1 | - | -5.9 | 0.62 | 0.67 | 1 | - | -6.9 |
| 1qbo | 4.48 | 1.32 | 1 | - | -6.6 | 3.65 | 0.78 | 1 | - | -7.7 |
| 1b8y | 0.95 | 1.11 | - | - | -9.5 | 1.40 | 0.67 | - | - | -10.1 |
| 1bwi | 5.40 | 1.14 | - | - | -5.4 | 6.14 | 0.55 | - | - | -6.2 |
| 1ciz | 2.01 | 0.45 | - | - | -9.7 | 1.39 | 1.19 | - | - | -10.8 |
| 1d8m | 1.03 | 1.11 | - | - | -7.4 | 1.37 | 1.49 | - | - | -8.0 |

^a rmsd (in Å): criterion of success of 2.0 Å. ^b rmsd (in Å): criterion of success: better than average rmsd; average rmsd between protein structures computed on the binding site residues: Factor Xa: 0.86 Å, trypsin: 0.90 Å, MMP-3: 0.92 Å. ^c Water molecules as proposed by FITTED; 1 and 0 define the presence or absence of each water molecule, respectively. ^d Score in arbitrary units. ^e The second best has an rmsd of 1.36 Å with a higher potential energy but a better score. ^f Poses with rmsd below 1.5 Å were found but given worse scores. Bold numbers highlight failures.

Cross Docking Study. In a real case study, medicinal chemists wish to design compounds de novo or to screen libraries of compounds that are not cocrystallized with the enzyme. Thus, a self-docking study is not representative of the real accuracy of docking programs. To properly evaluate the predictive power of FITTED, a set of cross-docking experiments was carried out next.

Each inhibitor was docked to the corresponding set of proteins in order to evaluate the impact of the protein conformation on the docking accuracy. First, each HIVP inhibitor of the monoprotonated subset was docked to each of the 5 protein structures, and the rmsds and scores were computed (Table 4). The data collected for the first five inhibitors revealed that the docking accuracy is greatly influenced by the protein conformation. The cross docking experiments carried out with the TK, FXa, and MMP inhibitors also showed a significant decrease in the accuracy relative to the self-docking study (Tables 5 and 6). In contrast, the other five HIVP and the trypsin inhibitors were accurately docked in most of the cases regardless of the protein structure used.

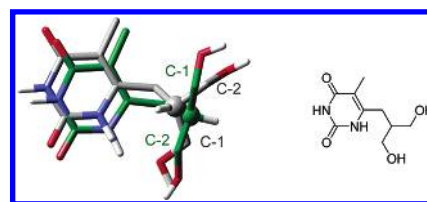


Figure 8. Docked (green) and crystal structure (gray) of 1e2p ligand; 2-D structure is shown to the right for clarity. Computed rmsd: 2.03 Å. The prochiral carbon is shown as a ball. Crystal structure (protein in gray, ligand in gray) and proposed docked model (protein in blue, inhibitor in green) for the 1d8m complex.

Overall, this cross-docking study confirms the need for a docking method that models the protein flexibility and/or the sensitivity of FITTED for the protein structure.

Docking to Multiple Conformations. The self-docking and cross-docking data can be used to simulate the docking to multiple conformations. The five (or nine for TK) final docked poses for each inhibitor (one per protein structure) were compared, and the best scoring pose is selected (Tables 4–6). In the case of the monoalcohols 1b6l, 1hpo, and 1pro, the self-docking led to the best score. The same observation was made with 5 out of the 9 investigated TK inhibitors and 4 out of the 5 FXa systems. However, as four of the five diols (1ajv, 1ajx, 1hvr, and 1hwr) were docked with good accuracy regardless of the protein structure, the prediction of the protein conformation was much poorer. Interestingly, although 1qbs was not accurately docked back to its corresponding protein structure, its correct binding mode associated with a better score was proposed when the protein structures 1ajx and 1hwr were employed.

Docking to Semiflexible and Flexible Proteins. Although the previous study intrinsically includes protein flexibility, it requires 5–9 experiments per compound and therefore implies the equivalent increase in required CPU time. The current version of FITTED offers to model the protein flexibility in a single experiment. Either one of the two options (docking to semiflexible and fully flexible proteins) described above can be selected. Adding the flexibility of the protein increases the complexity of the potential energy surface to explore, therefore making the conformational sampling more difficult. We therefore expected to observe a reduced accuracy when moving from rigid to flexible proteins.

In fact, 1hpo was misdocked to the semiflexible and flexible HIVP (Table 7), while an rmsd below 1.0 Å was recorded in the previous set of experiments. A close look at

the 10 runs revealed that the same misdocked pose was observed in 5 of the 10 proposed poses. It was also found that the experimentally observed pose had a worse score. These two observations ruled out the hypothesized bad convergence but pointed out a weakness in the scoring function. In contrast, 1qbs was misdocked to the rigid protein with the orientation of the seven-membered core reversed. As in the cross-docking study, a much better score and the right pose were predicted when the semiflexible and flexible proteins were used.

The computed protein rmsds between the superposed binding sites of the crystal structures 1of1 and 1e2k, 1e2p and 1e2k, and 1of1 and 1e2p of TK were 0.26, 0.59, and 0.59 Å, respectively. The computed rmsds for the other pairs of structures ranged from 0.80 to 1.16 Å with an average of 0.92 Å. When the semiflexible docking was used, the correct protein structure was picked among the possible nine in 4 cases (Table 8). When the 1of1 inhibitor was docked, the correct protein structure or a similar one was alternatively picked (5 runs each). In contrast, the protein structure was as good as the average when 1e2p was docked and worse than average when 1ki8 and 2ki5 were docked. Overall, the protein structure was predicted with an average rmsd of 0.44 Å for the eight successful dockings (lower than the average rmsd computed for each pair of structures). As discussed above, the Gln125 side chain of TK can adopt two distinct conformations. FITTED predicts the right conformation in 7 of the 8 successful docking cases. This is a good indicator of the predictive power of FITTED when the semiflexible option is selected. The docking to the fully flexible protein was less successful with an average rmsd of 0.78 Å but still below the average rmsd for the 9 protein structures (average rmsd = 0.92 Å). In this last study the Gln125 side chain was misoriented in 3 out of the 8 successful cases.

Data collected in Table 9 for FXa, trypsin, and MMP-3 shows that 1o2j, 1o3g, and 1o3i were properly docked, while 1f0r was misdocked to the semiflexible and flexible protein. Whether it was docked to the rigid, semiflexible, or flexible proteins, two alternative poses (rmsd ~1.0 Å or rmsd ~9.5 Å) were proposed for 1ezq. However, the wrong pose was assigned a better score when the fully flexible protein was used. The correct pose was much less observed (20% of the runs) than the wrong one. These results indicate that the global minimum may be located in a sharp and deep energy well of the potential energy surface that is difficult to find. In these series, the prediction of the occurrence of water molecules is good, while the protein structure prediction is more disappointing.

MMP-3 inhibitor 1bwi was consistently misdocked to the rigid, semiflexible, and fully flexible protein. More interestingly, 1d8m was misdocked to its corresponding protein crystal structure but properly docked to the 1b8y protein structure and to the semiflexible and fully flexible protein. A closer look at the 1d8m data showed that this inhibitor was properly docked when most dissimilar protein structures were used. This may indicate that some fine adjustments of the protein in the crystal structure of 1d8m would be required. In order to simulate these slight moves, FITTED selected a more appropriate protein structure. This may also indicate poor accuracy of the protein structure prediction for this enzyme or a poor description of the metal chelation.

Table 10. Docking Accuracy: Rigid Proteins

| | docking to protein ^a | docking to protein + water molecules ^b | | cross-docking |
|---------|---------------------------------|--|--------------------|---------------------|
| | ligand ^c | ligand ^c | water ^d | ligand ^c |
| success | 67% | 79% | 82% | 49% |

^a Water molecules removed prior to self-docking. ^b Bridging water molecules (see experimental section) were retained, and the function describing the interaction between ligand and water molecules was applied. ^c rmsd (in Å): criterion of success of 2.0 Å. ^d Criterion of success: occurrence predicted when ligand successfully docked.

This exhaustive docking study demonstrated that the scoring function can not only assign high scores to the experimentally observed pose but also discriminate between protein structures. It also shows that in specific cases such as 1qbs or 1d8m flexibility improves the accuracy over self-docking.

A comparison of the scores given to all docked poses showed that regardless of the protein flexibility method the score is within 1 unit for each compound. The scoring function is being further investigated, and improvements will be reported in due course.

DISCUSSION

Tables 10 and 11 summarize the accuracy observed throughout this study. It is worth mentioning that the tables show data for the top scoring poses only. This study was designed to assess the impact of the energy term used to model “displaceable” water molecules, on one hand, and the protein flexibility, on the other hand, on the accuracy of FITTED. First, a clear increase in accuracy was observed when the “displaceable” water molecules were added and validated the developed concept (Table 10). Overall, FITTED self-docked 79% of the inhibitors within 2.0 Å from the observed binding modes when the water was considered and only 67% when it was removed. In addition, the occurrence of water molecules was predicted with nearly 80% accuracy.

As a comparison, Kontoyianni and co-workers⁴ found GOLD and Glide as the most accurate programs with 69% and 58% of the compounds docked in a manner similar to the experimentally observed mode (referred to as “close” in Kontoyianni’s report), while LigandFit, FlexX, and DOCK showed poorer prediction powers. In another comparative study, Brooks and co-workers³ docked 73% and 46% of the compounds with rmsd below 2.0 Å with ICM and GOLD, respectively, while AutoDock, DOCK, and FlexX were less accurate. Rognan and co-workers² performed a similar study and found that Glide, GOLD, Surflex, and QXP docked 80–90% of the inhibitors within 2.0 Å from the observed pose, while FlexX, Fred, DOCK, and Slide showed lower accuracy (50–65%). Another study carried out by Perola and co-workers⁵ showed that Glide outperformed (61% within 2.0 Å) GOLD and ICM (48% and 45%, respectively). Although each study was based on a different set of protein–ligand complexes, our validation study shows that FITTED performed very well with accuracy as high or higher than the best performing docking programs. More importantly, FITTED allows for flexibility of the protein and displaceable water molecules to be accounted for, while GOLD includes water molecules but protein flexibility is restricted to the polar hydrogens, and Glide does not consider flexibility nor

Table 11. Docking Accuracy: Flexible Proteins

| | multiple conformations ^a | | | | semiflexible protein | | | | fully flexible protein | | | |
|----------------------|-------------------------------------|----------------------|-----|--------------------|----------------------|----------------------|-----|--------------------|------------------------|----------------------|-----|--------------------|
| | ligand ^b | protein ^c | | water ^d | ligand ^b | protein ^c | | water ^d | ligand ^b | protein ^c | | water ^d |
| success ^e | 79% | 47% | 76% | 73% | 73% | 27% | 61% | 82% | 73% | 0% | 73% | 81% |

^a Best scoring poses from self- and cross-docking studies (see text). ^b rmsd (in Å): criterion of success of 2.0 Å. ^c rmsd (in Å) calculated on successful dockings (ligand correctly docked): percentages of success given following two different criteria of success: exact protein structure (rmsd=0.0 Å), rmsd of the chosen protein below average rmsd of the protein pool. ^d Criterion of success: occurrence predicted. ^e The success rates are computed on the systems with the ligand successfully docked.

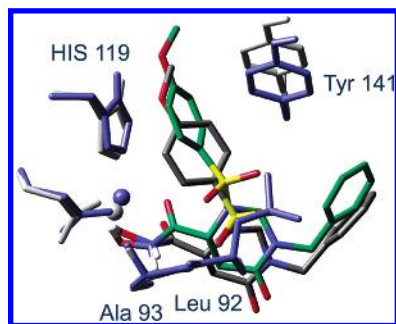


Figure 9. Crystal structure (protein in gray, ligand in gray) and proposed docked model (protein in blue, inhibitor in green) for the 1d8m complex.

water molecules. An exhaustive comparative study of the performance of FITTED with respect to other available docking programs is being undertaken in our research group and will be reported in due course.

With this data on hand, we turned our attention to the benefit and impact of the flexibility. Self-docking is the ideal case with the protein being molded to the ligand structure. In contrast, cross-docking attempts to combine ligands and protein structures that might not be perfectly matched, and it is well-known that docking programs perform poorer in this type of studies.¹⁰ In practice, docking experiments considering protein flexibility should be more accurate than cross-docking experiments and ideally as accurate as self-docking. In the present study, 49% and 79% accuracy were recorded for the cross- and self-docking experiments, respectively. Gratifyingly, the observed accuracy of FITTED when docking ligands to flexible proteins was close to that seen in self-docking. In addition, neither the prediction of the water molecule occurrence nor ligand pose was significantly affected by adding the protein flexibility as a slight drop in accuracy was observed when moving from rigid (self-docking) to flexible proteins.

A close look at the predicted protein structure revealed the good accuracy of our protocol. For instance, while 5–9 protein structures were used as input for each experiment, the correct conformation was selected for 9 of the 24 (37%) correctly docked systems when the semiflexible protein was used. If a random selection was used, then 11–20% of the cases would present the correct structure. When considering the successful docking experiments (ligand pose accurately predicted), average protein rmsds of 0.50 Å and 0.78 Å were recorded when using the semiflexible or fully flexible proteins respectively.

Unexpected results were also recorded. In two cases (MMP-3 1d8m and HIVP 1qbs), docking to flexible proteins was more accurate than self-docking. Figure 9 shows a superposition of the crystal structure of 1d8m and the docked structure when the semiflexible option was used. When the

rigid protein was used, a wrong binding mode was proposed, while a better prediction was observed with the flexible protein. The accurately docked pose is slightly translated relative to the crystal structure and indicates weak hydrogen bonds but strong hydrophobic/ π -stacking interactions with His119 and Tyr141 side chains located in the S₁' pocket. In order to induce this interaction pattern, the S₁' pocket must be more closed in the modeled structure than in the crystal structure to encompass the ligand better. This discrepancy may show that the hydrogen-bonding contribution to the binding is underestimated or that the van der Waals interactions are overestimated.

In the case of 1qbs, no clear explanation was found. A close look at the docked and crystal complexes does not reveal any specific movement or steric clash. We therefore believe that the potential energy function may find some nuances that can be detrimental to the correct pose. Again, slight adjustments of the crystal structures may be necessary prior to or upon docking. This hypothesis claims that the crystal structure of 1qbs included some discrepancies that induced some slight van der Waals repulsions, preventing a good score when docking 1qbs inhibitors to the 1qbs protein structure. These repulsions would vanish when the flexible protein was used and much lower (better) scores were recorded. Possible strategies to address this issue are the use of relaxed structures (as proposed in Glide³⁴), soft structures as proposed by Shoichet and co-workers,³⁵ or flexible structures as shown in this study. These two unexpected results may reveal some inaccuracies of the scoring function.

Overall, this study revealed the accuracy of FITTED to dock inhibitors to flexible and partially solvated proteins and validated it with this set of representative protein–inhibitor complexes.

CONCLUSION

We have developed FITTED 1.0, a unique docking program that accounts for both protein flexibility and bridging water molecules. The flexibility is handled by a genetic algorithm based on various genetic operators specific to FITTED (i.e., designed crossover operator, focused mutation, filters). Modifications to the initial genetic algorithm have been made to increase the speed and accuracy by orienting the docking toward “favored” poses (e.g., poses within the cavity and fulfilling constraints). We have also implemented a new potential energy term that accurately accounts for dynamically bound water molecules. Application of FITTED to the docking of a variety of protein/inhibitors complexes resulted in proposed docked poses within 2.0 Å from the observed binding modes in 73–76% of the cases using flexible or rigid proteins, respectively. The accurate prediction of the occurrence and the need for displaceable

water molecules was also demonstrated. Finally, the protein structures were predicted with reasonable accuracy.

Our initial studies led to a method that docked each compound within 0.5–20 h when not considering rotation and orientation of the ligand as part of the chromosomes. FITTED, which now explores the entire conformational space of the ligands, considers protein flexibility and displaceable water molecules and docks all the tested compounds within 3 h on a similar processor. Further studies are in progress to reduce by a factor of 10 or more the required CPU time which is still not appropriate for virtual screening and to improve the scoring function.

EXPERIMENTAL SECTION

Preparation of the Training Set. General Considerations. The protein/ligand complexes were retrieved from the Protein Data Bank³⁶ or from the PDBbind database.^{37,38} The complexes were selected for the occurrence of water molecules, for the flexibility of the protein structure, the diversity of the ligands, resolutions lower than 2.5 Å, and good binding affinities. At least 4 structures for each system (MMP, HIVP, TK, FXa, trypsin) were looked for. The complexes were set up using Maestro (v9.0) and/or InsightII (v2005). The set of complexes from the same family was superimposed prior to its use with FITTED. In order to be able to use FITTED with more than one protein structure, the sequences have to be identical. Therefore, some minor mutations (often far from the binding site) were achieved (e.g., Arg14 into Lys14 in HIVP complex 1b6l), missing side chains were reconstructed (e.g., Arg220 in 1qhi), and names of residues were made identical (e.g., Glu124A into Glu124 in 1nfu). Hydrogens were next added and optimized by energy minimization. All the nonconserved waters were removed, and missing key water molecules were added by analogy with other structures when applicable. For instance, the water 301 observed in many HIVP/inhibitor complexes is displaced in 1ajv and was added to the 1ajv protein structure, and its position was optimized by energy minimization using AMBER94 as a force field. The naming of the water molecules is made homogeneous within each set. Each protein is then saved as a mol2 file and processed using ProCESS to assign protein atom types and charges. Each ligand was charged using Sybyl Gasteiger–Hückel charges and processed using SMART. Large grids of spheres were prepared as well as constraint files. These constraints were loose in order to orient and speed up the docking (as previously described) but not bias the results. Diameters of the constraint spheres as large as 8.0 Å were used.

HIV-1 Protease Inhibitor/Protein Complexes. HIVP complexes following the criteria defined above were retrieved from the PDB. 1eby (crystal structure resolution: 2.29 Å), 1hpo (2.50 Å), 1hpu (1.90 Å), and 1pro (1.80 Å) were superimposed onto 1b6l (1.75 Å), and one catalytic aspartic acid side chain was protonated. 1ajx (2.00 Å), 1hvr (1.80 Å), 1hwr (1.80 Å), and 1qbs (1.80 Å) were superimposed onto 1ajv (2.00 Å), and the two catalytic aspartic acid side chains were protonated following the experimental study of similar complexes.³³ A water molecule hydrogen-bonding to both Ile50 NHs was kept when present in the crystal structures and added by analogy when missing. The constraint applied imposes a polar group to be located close to

the catalytic site. As some of the inhibitors have a large number of rotatable bonds, initial populations of 200 individuals were used in all 10 cases.

Thymidine Kinase Inhibitor/Protein Complexes. All the available TK inhibitor/protein complexes were retrieved from the PDB and filtered. A final set of nine structures was used: 1e2k (1.70 Å), 1e2p (2.50 Å), 1ki3 (2.37 Å), 1ki4 (2.34 Å), 1ki7 (2.20 Å), 1ki8 (2.20 Å), 2ki5 (1.90 Å), 1of1 (1.95 Å), and 1qhi (1.90 Å). Six key water molecules were considered as discussed in the text. The constraint imposes polar groups to be located within the catalytic site.

Factor Xa Inhibitor/Protein Complexes. These complexes were retrieved from the PDBbind database. 1f0r (2.10 Å), 1fjs (1.92 Å), 1nfu (2.05 Å), and 1xka (2.30 Å) were superimposed onto 1ezq (2.20 Å). In this set, two key water molecules were identified and added when missing to the protein structures. The constraint imposes a polar group to be located close to the Asp189 side chain. Initial populations of 200 individuals were used.

Trypsin Inhibitor/Protein Complexes. 1o2j (1.65 Å), 1o3g (1.55 Å), 1o3i (1.51 Å), and 1qbo (1.80 Å), were superimposed onto 1f0u (1.90 Å). In this case a single water molecule interacting with Leu227 was considered. The constraint imposes a polar group to be located close to the Asp171 side chain. Initial populations of 200 individuals were used.

MMP Inhibitor/Protein Complexes. 1bwi (1.80 Å), 1ciz (1.64 Å), and 1d8m (2.44 Å) were superposed onto 1b8y (2.00 Å). No water molecules were retained. The constraint imposes a polar group to be located close to the catalytic zinc atom. Specific zinc (van der Waals and metal chelation) and hydroxamic acid (internal energy) parameters were added to the force field. The LJ12-10 potential parameters used for the zinc atom were designed to reproduce the observed energy of zinc chelation.³⁹

Docking Study. Self-Docking, Semiflexible Protein, and Fully Flexible Protein. In the first of these three sets of runs (self- and cross-docking, docking to multiple conformations), one single protein structure was used as an input to evaluate the accuracy of the docking algorithm. In the second set (docking to semiflexible proteins), the protein structure was restricted to 4-9 input conformations. In the third set (docking to fully flexible protein), the protein structures were composite of 4-9 input conformations. A typical keyword file with all the default parameters is given as Supporting Information. The default parameters (e.g., 10 runs, population size of 100 individuals) were used unless otherwise stated.

ProCESS and FITTED Parameters. The ensemble of spheres cavity of the binding site were centered on the center of the cavity and did not exceed 28 Å long. The grid resolution was 1.5 Å.

ACKNOWLEDGMENT

We thank Virochem Pharma for financial support and a scholarship to C.R.C. as well as the Canadian Foundation for Innovation for financial support through the New Opportunities Fund program. P.E. is supported by a scholarship from Canadian Institutes of Health Research (Strategic Training Initiative in Chemical Biology). We thank RQCHP for generous allocation of computer resources.

Supporting Information Available: Typical keyword files for FITTED and ProCESS and a detailed description of the validation set (PDB codes, structures, K_i 's). This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- (1) Rester, U. Dock around the Clock — Current Status of Small Molecule Docking and Scoring. *QSAR Comb. Sci.* **2006**, 25, 605–615.
- (2) Bissantz, C.; Folkers, G.; Rognan, D. Protein-Based Virtual Screening of Chemical Databases. 1. Evaluation of Different Docking/Scoring Combinations. *J. Med. Chem.* **2000**, 43, 4759–4767.
- (3) Bursulaya, B. D.; Totrov, M.; Abagyan, R.; Brooks, C. L., III. Comparative study of several algorithms for flexible ligand docking. *J. Comput.-Aided Mol. Des.* **2003**, 17, 755–763.
- (4) Kontoyianni, M.; McClellan, L. M.; Sokol, G. S. Evaluation of Docking Performance: Comparative Data on Docking Algorithms. *J. Med. Chem.* **2004**, 47, 558–565.
- (5) Perola, E.; Walters, W. P.; Charifson, P. S. A Detailed Comparison of Current Docking and Scoring Methods on Systems of Pharmaceutical Relevance. *Proteins: Struct., Funct., Bioinf.* **2004**, 56, 235–249.
- (6) Kellenberger, E.; Rodrigo, J.; Muller, P.; Rognan, D. Comparative Evaluation of Eight Docking Tools for Docking and Virtual Screening Accuracy. *Proteins: Struct., Funct., Bioinf.* **2004**, 57, 225–242.
- (7) Cummings, M. D.; DesJarlais, R. L.; Gibbs, A. C.; Mohan, V.; Jaeger, E. P. Comparison of Automated Docking Programs as Virtual Screening Tools. *J. Med. Chem.* **2005**, 48, 962–976.
- (8) Kontoyianni, M.; Sokol, G. S.; McClellan, L. M. Evaluation of Library Ranking Efficacy in Virtual Screening. *J. Comput. Chem.* **2005**, 26, 11–22.
- (9) Cavasotto, C. N.; Abagyan, R. A. Protein Flexibility in Ligand Docking and Virtual Screening to Protein Kinases. *J. Mol. Biol.* **2004**, 337, 209–225.
- (10) Osterberg, F.; Morris, G. M.; Sanner, M. F.; Olson, A. J.; Goodsell, D. S. Automated Docking to Multiple Target Structures: Incorporation of Protein Mobility and Structural Water Heterogeneity in AutoDock. *Proteins: Struct., Funct., Genet.* **2002**, 46, 34–40.
- (11) Murray, C. W.; Baxter, C. A.; Frenkel, A. D. The Sensitivity of the Results of Molecular Docking to Induced Fit Effects: Application to Thrombin, Thermolysin and Neuraminidase. *J. Comput.-Aided Mol. Des.* **1999**, 13, 547–562.
- (12) Murray, C. W.; Baxter, C. A.; Frenkel, A. D. The Sensitivity of the Results of Molecular Docking to Induced Fit Effects: Application to Thrombin, Thermolysin and Neuraminidase. *J. Comput.-Aided Mol. Des.* **1999**, 13, 547–562.
- (13) Erickson, J. A.; Jalaie, M.; Robertson, D. H.; Lewis, R. A.; Vieth, M. Lessons in Molecular Recognition: the Effects of Ligand and Protein Flexibility on Molecular Docking Accuracy. *J. Med. Chem.* **2004**, 47, 45–55.
- (14) Carlson, H. A. Protein Flexibility and Drug Design: How to Hit a Moving Target. *Curr. Opin. Chem. Biol.* **2002**, 6, 447–452.
- (15) Schnecke, V.; Kuhn, L. A. Virtual Screening with Solvation and Ligand-Induced Complementarity. *Perspect. Drug Discovery Des.* **2000**, 20, 171–190.
- (16) Claussen, H.; Buning, C.; Rarey, M.; Lengauer, T. FlexE: Efficient Molecular Docking Considering Protein Structure Variations. *J. Mol. Biol.* **2001**, 308, 377–395.
- (17) Zavodszky, M. I.; Lei, M.; Thorpe, M. F.; Day, A. R.; Kuhn, L. A. Modeling Correlated Main-Chain Motions in Proteins for Flexible Recognition. *Proteins: Struct., Funct., Bioinf.* **2004**, 57, 243–261.
- (18) While this manuscript was in preparation, the beta version of Autodock4 came out (<http://autodock.scripps.edu/doc/AutoDock4>). Autodock4 also models the side-chain flexibility.
- (19) Sherman, W.; Day, T.; Jacobson, M. P.; Friesner, R. A.; Farid, R. Novel Procedure for Modeling Ligand/Receptor Induced Fit Effects. *J. Med. Chem.* **2006**, 49, 534–553.
- (20) Moitessier, N.; Westhof, E.; Hanessian, S. Docking of Aminoglycosides to Hydrated and Flexible RNA. *J. Med. Chem.* **2006**, 49, 1023–1033.
- (21) Moitessier, N.; Henry, C.; Maigret, B.; Chapeleur, Y. Combining Pharmacophore Search, Automated Docking, and Molecular Dynamics Simulations as a Novel Strategy for Flexible Docking. Proof of Concept: Docking of Arginine-Glycine-Aspartic Acid-like Compounds into the $\alpha_1\beta_3$ Binding Site. *J. Med. Chem.* **2004**, 47, 4178–4187.
- (22) Moitessier, N.; Therrien, E.; Hanessian, S. A Method for Induced-fit Docking, Scoring and Ranking of Flexible ligands. Application to Peptidic and Pseudopeptidic BACE 1 Inhibitors. *J. Med. Chem.* **2006**, 49, 5885–5894.
- (23) Wei, B. Q.; Weaver, L. H.; Ferrari, A. M.; Matthews, B. W.; Shoichet, B. K. Testing a Flexible-Receptor Docking Algorithm in a Model Binding Site. *J. Mol. Biol.* **2004**, 337, 1161–1182.
- (24) *CDISCOVER*, 98.0; Accelrys, Inc.: San Diego, CA, 2001.
- (25) Fletcher, R.; Reeves, C. M. Function Minimization by Conjugate Gradients. *Comp. J.* **1964**, 7, 149–154.
- (26) Morris, G. M.; Goodsell, D. S.; Halliday, R. S.; Huey, R.; Hart, W. E.; Blelew, R. K.; Olson, A. J. Automated Docking Using a Lamarckian Genetic Algorithm and an Empirical Binding Free Energy Function. *J. Comput. Chem.* **1998**, 19, 1639–1662.
- (27) http://www.rcsb.org/pdb/file_formats/pdb/pdbguide2.2/guide2.2_frame.html (accessed Dec 2006).
- (28) Corey, E. J.; Wipke, W. T. Computer-Assisted Design of Complex Organic Syntheses. *Science* **1969**, 166, 178–192.
- (29) Gasteiger, J.; Marsili, M. Iterative partial equalization of orbital electronegativity—a rapid access to atomic charges. *Tetrahedron* **1980**, 36, 3219–3228.
- (30) Haupt, R. L. Optimization of aperiodic conducting grids. *11th Annual Rev. Progress in Applied Computational Electromagnetics Conference*, Monterey, CA, 1995.
- (31) Weiner, S. J.; Kollman, P. A.; Nguyen, D. T.; Case, D. A. An All Atom Force Field for Simulations of Proteins and Nucleic Acids. *J. Comput. Chem.* **1986**, 7, 230–252.
- (32) See, for example: Verdonk, M. L.; Chessari, G.; Cole, J. C.; Hartshorn, M. J.; Murray, C. W.; Nissink, J. W. M.; Taylor, R. D.; Taylor, R. Modeling Water Molecules in Protein-Ligand Docking Using GOLD. *J. Med. Chem.* **2005**, 48, 6504–6515.
- (33) Yamazaki, T.; Nicholson, L. K.; Torchia, D. A.; Wingfield, P.; Stahl, S. J.; Kaufman, J. D.; Eyermann, C. J.; Hedge, C. N.; Lam, P. Y. S.; Ru, Y.; Jadhav, P. K.; Chang, C.-H.; Weber, P. C. NMR and X-ray Evidence That the HIV Protease Catalytic Aspartyl Groups Are Protonated in the Complex Formed by the Protease and a Non-peptide Cyclic Urea Inhibitor. *J. Am. Chem. Soc.* **1994**, 116, 10791–10792.
- (34) Friesner, R. A.; Banks, J. L.; Murphy, R. B.; Halgren, T. A.; Klicic, J. J.; Mainz, D. T.; Repasky, M. P.; Knoll, E. H.; Shelley, M.; Perry, J. K.; Shaw, D. E.; Francis, P.; Shenkin, P. S. Glide: A New Approach for Rapid, Accurate Docking and Scoring. 1. Method and Assessment of Docking Accuracy. *J. Med. Chem.* **2004**, 47, 1739–1749.
- (35) Ferrari, A. M.; Wei, B. Q.; Costantino, L.; Shoichet, B. K. Soft Docking and Multiple Receptor Conformations in Virtual Screening. *J. Med. Chem.* **2004**, 47, 5076–5084.
- (36) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Bourne, P. E.; Shindyalov, I. N. The Protein Data Bank. *Nucleic Acids Res.* **2000**, 28, 235–242.
- (37) Wang, R.; Fang, X.; Lu, Y.; Wang, S. The PDBbind database: Collection of binding affinities for protein-ligand complexes with known three-dimensional structures. *J. Med. Chem.* **2004**, 47, 2977–2980.
- (38) Wang, R.; Fang, X.; Lu, Y.; Yang, C. Y.; Wang, S. The PDBbind database: Methodologies and updates. *J. Med. Chem.* **2005**, 48, 4111–4119.
- (39) Tiraboschi, G.; Greshh, N.; Giessner-Pretre, C.; Pedersen, L. G.; Deerfield, D. W. Parallel Ab Initio and Molecular Mechanics Investigation of Polycoordinated Zn(II) Complexes with Model Hard and Soft Ligands: Variations of Binding Energy and of Its Components with Number and Charges of Ligands. *J. Comput. Chem.* **2000**, 21, 1011–1039.

CI6002637