IEEE/ACM TRANSACTIONS ON COMPUTATIONAL BIOLOGY AND BIOINFORMATICS

CaverDock: A Novel Method for the Fast Analysis of Ligand Transport

Jiří Filipovič, Ondřej Vávra, Jan Plhák, David Bednář, Sérgio M. Marques, Jan Brezovský, Luděk Matyska, Jiří Damborský

Abstract—Here we present a novel method for the analysis of transport processes in proteins and its implementation called CaverDock. Our method is based on a modified molecular docking algorithm. It iteratively places the ligand along the access tunnel in such a way that the ligand movement is contiguous and the energy is minimized. The result of CaverDock calculation is a ligand trajectory and an energy profile of transport process. CaverDock uses the modified docking program Autodock Vina for molecular docking and implements a parallel heuristic algorithm for searching the space of possible trajectories. Our method lies in between the geometrical approaches and molecular dynamics simulations. Contrary to the geometrical methods, it provides an evaluation of chemical forces. However, it is far less computationally demanding and easier to set up compared to molecular dynamics simulations. CaverDock will find a broad use in the fields of computational enzymology, drug design and protein engineering. The software is available free of charge to the academic users at https://loschmidt.chemi.muni.cz/caverdock/.

Index Terms—molecular docking, tunnel analysis, ligand transport, drug design, numerical optimization, restrained force field, volume discretization

1 INTRODUCTION

Understanding protein-ligand interactions is of great importance in many fundamental biochemical processes as well as in various applications. For example, to study a ligand that may inhibit protein function allowing a virus to attack a cell or to design inhibitors blocking tunnels and channels as a new paradigm in drug design [1]. The ligand interacts with protein in its active or binding site - the functional site of a protein. Simulation of ligand binding (entering the active site and forming a stable complex) and unbinding (release of a ligand from a stable complex) helps in many practical applications. It allows to search for ligands which are more likely to bind to a particular protein; modify a ligand to bind faster or with higher affinity or modify the protein to ease or disallow the ligand binding. Many proteins have binding sites buried inside their cores, which implies that a ligand must traverse through a *tunnel* or a *channel*¹ before it can bind to the functional site in the protein. In such cases, we need to analyze whether the ligand is likely to pass through the tunnel or channel into the protein core.

All chemical systems, such as the proteins interacting with ligands, follow the second law of thermodynamics: they tend to minimize their potential energy. In practice, the most probable *conformation* of the molecules (i. e. spatial

- O. Vávra, D. Bednář, S. M. Marques, J. Brezovský, J. Damborský are with Department of Experimental Biology and RECETOX, Faculty of Science, Masaryk University. E-mail: {393643, 222755, smarques, brezovsky}@mail.muni.cz, jiri@chemi.muni.cz.
- J. Filipovič and O. Vávra contributed equally to this work.

1. In the following text, we will speak about tunnels for simplicity. However, the mechanisms of a ligand passing through a channel are the same as for a tunnel. position of their atoms) is the one with the lowest potential energy. However, it is also possible for molecules to make a transition from one local minimum to another, depending on system temperature and height of energetic barrier – the smaller energetic barrier, the more probable is the transition. In the molecular modeling methods, the function estimating the potential energy for a given conformation is called a *force field*. The analysis of the potential energy given by the force field allows us to compute the probability of some conformation to appear in the real-world chemical system. It allows to predict whether or how fast some chemical process (such as a ligand passing through a tunnel) can occur at a given temperature.

To study the ligand binding or unbinding, we need to evaluate the potential energy of the ligand passing from protein surface through the tunnel into the active site or vice versa. The ligand binds in an active site if there is a strong local energetic minimum and it passes through the tunnel if there is no significant energy barrier along the way (the gradient of potential energy is more or less decreasing from tunnel entrance to its binding site). When a tunnel contains some strong repulsive barrier, the ligand is likely not to pass through the tunnel. Note that the energetic profile of the tunnel is unique with respect to the ligand, as it reflects the specific ligand-protein interactions occurring during the ligand passage.

The ligand binding to a protein's active site or binding site is usually computed by *molecular docking*. A molecular docking algorithm traverses the conformation space of the protein-ligand complex and searches for energetic minima [2], [3], [4], [5]. The result of the molecular docking is the structure of the protein-ligand complexes together with an estimation of the respective free energy of binding. Thus, the users can learn which ligand binds with the lowest energy or study the orientation of a ligand in a protein active

J. Filipovič, J. Plhák and L. Matyska are with the Institute of Computer Science, Masaryk University, Botanická 68a, Brno, Czech Republic. Email: {fila, 408420, ludek}@mail.muni.cz

2

site. However, the molecular docking computes only the lowest energy positions of a ligand within some region of a protein, and thus it is not suitable for the study of the ligand transport through a protein tunnel.

In this paper, we present a novel method for computation of the binding free energy variation and the ligand *trajectory* (movement of the ligand's atoms along the tunnel). It allows to study the processes of ligand binding and unbinding through the tunnels of any protein. Our method is based on a molecular docking algorithm – it iteratively docks the ligand along the previously calculated tunnel and at each point it evaluates its binding free energy. Our docking works with restraints. It uses a combination of the chemical force-field from AutoDock Vina [2] to compute the binding free energy of the protein-ligand complex and a newly developed restraints, which restricts the ligand positions to specified part of the tunnel and to vicinity of defined conformation. Note that restraints produce penalization energy, which is used to keep a ligand in a suitable position, but the binding free energy reported by CaverDock is based on the original AutoDock Vina force field only. With the restraints, a contiguous ligand trajectory with arbitrary step size (the maximal change in ligand's atoms position between two consecutive conformations) can be generated. Thus, the position of the ligand within the tunnel can be constrained to a defined area. Since there may exist many possible paths through a given tunnel, the paths are searched using a heuristic algorithm with backtracking. Our method is implemented in the user-friendly software tool CaverDock, which uses parallel architecture to maximize the performance of the ligand transit computation (from minutes to a few hours using a desktop computer).

This paper focuses primarily on the computer-science topics: it introduces our method and its implementation. It also presents basic evaluation, which illustrates that the produced trajectories and energetic profiles can be obtained in a reasonable time. The paper targeting the CaverDock user community, focused on the biochemical topics (setting the calculation and interpretation of results, evaluation and benchmarking CaverDock with realistic use cases on many protein-ligand pairs) is prepared in parallel with this paper [6], [7].

The rest of the paper is organized as follows. Section 2 summarizes the work related to our paper and describes the difference between our method and state-of-the-art. The high-level overview of our method is given in Section 3. Following three sections discuss the method in detail: Section 4 introduces the algorithm for tunnel discretization, Section 5 describes our modifications of docking algorithm using restraints for ligand position and Section 6 discuss the searching of trajectory space. The evaluation of our implementation is given in Section 7. We conclude and sketch future work in Section 8.

2 RELATED WORK

A very fast approximation of biomolecules represents the atoms as solid macro-world objects. The transport process of a ligand through a tunnel in a protein is then studied analogously to a macro-world objects' mechanics: the molecular shape is formed by spheres representing the atoms (which may be connected through flexible joints), neglecting any chemical forces, such as electrostatics, hydrogen bonds or solvation effects.

The majority of the geometry-based approaches analyze the tunnel only, without generating a ligand trajectory [8], [9], [10]. Those software tools take a protein or multiple conformations of the protein as the input and generate the geometry of the tunnel. The ability to transport a ligand is then judged based on tunnel geometry. A comprehensive study of different geometrical methods can be found in [11], [12].

A different approach to the geometric analysis is taken in MoMA-LigPath [13]. The ligand transport is studied here by an algorithm inspired by robotic motion planning. The protein and ligand are understood as mechanical objects, which are partially flexible as they may change the dihedral angles. The algorithm searches for a ligand trajectory from the active site to the tunnel entrance by moving the ligand and the flexible parts of the receptor. Such algorithm allows to detect parts of the receptor which need to be moved to allow the ligand to pass through the tunnel. However, it does not use a chemical force field, so there is no quantitative information showing how difficult is for the ligand to pass the tunnel due to chemical interactions (attractions and repulsions). Moreover, the induced movement of the ligand and the flexible parts may be unrealistic, as with the chemical forces different movements may be preferred.

The molecular dynamics (MD) uses an empirical force field to model the physical properties of the atoms and their interactions in time. There are many well-established software tools for MD, such as Amber [14] or Gromacs [15]. However, it is not practical to model the transportation of a ligand through a tunnel with classical MD, as the simulation time is often extremely long.Therefore, various modifications of MD are used to speed-up the process.

The metadynamics is an enhanced sampling technique which introduces biases in the form of repulsion energy on the already visited parts of the conformational space, such as the conformations of a molecule or the positions of a ligand within a tunnel [16]. The bias is computed according to the simulation state defined in term of collective variables (a small number of variables describing the simulation space). The metadynamics can be used to pass a ligand through a tunnel in a protein [17] and evaluate the thermodynamics and kinetics of the process. It explores simulation states much faster than MD, however, comparing to the geometrical approaches, it is still much more computationally demanding. Moreover, an expert user has to setup the metadynamics computation properly, as an incorrect definition of the collective variables may lead to inefficient biassing.

Another technique based on MD allowing to simulate transportation through a tunnel is the steered MD [18]. With steered-MD, the external force is applied to a ligand such that it is pulled from or to the tunnel. The technique is, similarly to metadynamics, more computationally demanding than the geometrical methods. An expert user has to set up how the external force is applied, otherwise, there can be a false bottleneck observed (e.g. when a ligand is pulled in the wrong direction against the protein backbone).

The molecular docking has been developed for eval-

FILIPOVIČ, VÁVRA et al.: CAVERDOCK: A NOVEL METHOD FOR THE FAST ANALYSIS OF LIGAND TRANSPORT

uation of the ligand binding free energy in the protein active site. It performs the search of many protein-ligand conformations and returns the ones in significant energetic minima. Thus, it is possible to study if the ligand binds preferably in the active site, or compare the minima of different ligands (i.e. to identify which ligands are more likely to interact with the protein). Many docking software tools are well-established and widely-used in the scientific community [2], [3], [4], [5]. Molecular docking is not suitable to be directly used for analysis of ligand transport, as it samples conformation space coarsely to find significant minima, but does not search the ligand path into the minima. However, it can be used as a basis for docking-based tools.

Similarly to our tool, molecular docking is used to analyze the transport process in SLITHER [19]. However, the SLITHER does not employ restrained docking. Instead, the movement of the ligand is induced by the biasing force at the already visited positions. Therefore, there is no mechanism to ensure the ligand movement is contiguous or at least provides fine-grained sampling of the trajectory in the tunnel – it may jump over bottlenecks without sampling the energy barriers or even jump into a different tunnel. Moreover, there is no sophisticated tunnel geometry analysis and the ligand is moved along the *y*-axis only. Therefore, SLITHER cannot be reliably used for highly curved tunnels, e.g. U-shaped.

3 METHOD OVERVIEW

In this section, we describe the basic concept of our method. The more detailed discussion will be given in the following sections. The method is based on a driven step-by-step movement of the ligand through the tunnel.We first discretize the tunnel into a set of discs, so the ligand movement through the tunnel can be driven, i. e. it is possible to define a ligand position in the tunnel and thus also movement "forward" and "backward" in the tunnel. After the discretization, the ligand is iteratively docked into consecutive positions along the tunnel, allowing to compute binding or unbinding trajectory.

3.1 Tunnel Discretization

To drive the ligand movement in the tunnel, we need to restrict the space where the ligand can be placed in each step. We use the tunnel geometry approximated by a sequence of spheres as the input. Such sequence can be obtained from Caver [8] or a similar tool. The sequence of spheres is then transformed into a sequence of *n* disks $\theta_1, \ldots, \theta_n$. We create the disks by cutting the tunnel into slices of an upper-bound thickness. The path of the ligand through the tunnel can be defined as the iterative placement of the selected ligand's atom to consecutive disks. Note that an arbitrary atom of the ligand can be selected, but it must be the same for the whole tunnel trajectory.

3.2 Docking with Restraints

The ligand conformation λ is defined by the Cartesian position of its atoms: $\lambda = \{a_i\}_{i=1}^m$. Having a discretization



3

Fig. 1. Schematic 2D view of traversing tunnel, where the selected ligand's atom is placed onto consecutive disks. As no contiguous movement of the ligand is required, the ligand flips between disks θ_6 and θ_7 , thus the small geometrical bottleneck between those disks is not detected.



Fig. 2. Schematic 2D view of a ligand traversing a tunnel. The ligand is depicted in black, its previous position used as a pattern is shown in grey. Restricting the movement of atoms causes the geometrical bottleneck between θ_6 and θ_7 to be detected when the ligand passes from θ_7 to θ_8 , as can be seen in the last figure.

of our tunnel, we can select an atom of the ligand $a_c \in \lambda$, which is placed onto any position of the selected disk θ :

$$a_c \in \theta$$
 (1)

We say the ligand is docked onto the disc when its atom a_c lies onto the disc. By placing the atom a_c onto consecutive discs $\theta_1, \ldots, \theta_n$, we force the ligand to move through the tunnel. Such ligand trajectory samples the tunnel without large gaps (i.e. the ligand cannot overcome very narrow bottlenecks or even jump to different tunnel), however, the trajectory is not contiguous (the ligand can e.g. rotate freely). We use this non-contiguous trajectory to compute the *lower-bound* energy profile of the ligand transport. The example of such trajectory is depicted in Figure 1. As we can see, atom a_c is stuck to the disk and by moving through the tunnel, the sampling of the transport process is obtained. However, the ligand may perform non-contiguous movement: it flips between disc θ_6 and θ_7 .

The contiguous trajectory can be computed by restricting the movement of each atom by constant δ . When a new ligand conformation λ_{i+1} is generated, the distance of each atom from its previous position in λ_i is upper-bound:

$$\forall j \in [1,m] : |a_j - b_j| < \delta \tag{2}$$

where $a_j \in \lambda_i, b_j \in \lambda_{i+1}$. We say λ_i is the pattern restraining the position of λ_{i+1} , formally: $\lambda_{i+1} \in \Delta \lambda_i$, when Eq. 2 holds.

IEEE/ACM TRANSACTIONS ON COMPUTATIONAL BIOLOGY AND BIOINFORMATICS



Fig. 3. Schematic 2D view of a ligand traversing a tunnel, where the ligand is stuck at a bottleneck (left figure). When a different orientation of the ligand is selected, the ligand may pass without reaching the barrier (right figure).

We can use the pattern restraint when a new position of the ligand is generated. Let λ^i represents the ligand conformation docked onto disc θ_i . When the ligand position λ^{i+1} on the disc θ_{i+1} is searched, the pattern restraint ensures that $\lambda^{i+1} \in \Delta \lambda^i$, thus, transition between discs is contiguous (upper-bound by δ). Note that the discs must be generated such that the distance between the discs must be lower than δ . The example of using pattern restraint is depicted in Figure 2. As we can see, the pattern disallows the ligand to flip (as exemplified by the movement from disc 6 to disc 7 in Figure 1), and the small geometrical bottleneck is detected.

3.3 Trajectory Search

4

The contiguous trajectory can be obtained by iterative docking onto the disks with restricted changes in the position of all atoms by a pattern restraint. However, we want to allow the ligand to optimize its position at each disc to find a local energetic minimum. This minimum may be unreachable after one step when the ligand movement is restricted by a pattern. Thus, we search for the ligand trajectory, where multiple conformations may be docked onto the same disc. More precisely, having the conformation $\lambda_{j+1}^i \in \Delta \lambda_j^i, \lambda_{j+2}^i \in \Delta \lambda_{j+1}^i, \ldots$ until the energy of the new conformations is improved. We call these steps the optimization steps, as they allow the ligand to find a low-energy position on the disc, which may not be feasible immediately after the transition from θ_{i-1} to θ_i .

The ligand movement described above prefers the transition where the ligand follows the strongest energy gradient locally between following steps. Although this scenario is the most probable in real-world systems, the ligand may occasionally make a transition to some different conformation, which may allow it to pass the energy barrier with lower energy. Consider the case depicted in Figure 3. Depending on its orientation, the ligand may or may not get through the tunnel bottleneck. Thus we need to search multiple variants of the ligand trajectory.

The number of contiguous trajectories may be very high – the transition to a new disk may change the ligand position, orientation, and conformation (relative position of the atoms within the ligand). The exhaustive search of possible trajectories is not feasible due to the time required to dock ligands with restraints (typically hundreds of milliseconds). Thus, we have introduced a simple heuristic. We move the ligand only in one direction in the tunnel (e.g. from the binding site to the protein surface). When the binding free energy of λ_j^i is significantly higher than the binding free energy of some known conformation λ_{low}^i (i.e. ob-

tained during lower-bound trajectory computation), we set $\lambda_j^i = \lambda_{low}^i$, and search the trajectory moving the ligand backwards to previous disks $\theta_{i-1}, \theta_{i-2}, \ldots$. The backtracking ends after the forward and backward trajectories converge, or after the beginning of the tunnel is reached. Note that the resulting trajectory still follows only one direction. When the backtrack is used the trajectory is reversed and integrated into a forward trajectory.

The situation when the backtracking trajectory converges with a forward trajectory (i.e. $\lambda_{backtrack}^{i} \in \Delta \lambda_{forward}^{i}$) allows to join both trajectories. The optimization of the ligand position moves it to a minimum at the current disc which allows convergence in many cases. However, the ligand may need to overcome some energetic barrier to converge. Therefore, we use also an explicit convergence process: a weak force is applied to the ligand in the backtrack trajectory in order to pull its position to the vicinity of the ligand in the forward trajectory.

3.4 CaverDock Workflow

From the user's perspective, CaverDock is a command-line tool taking the molecules' structures and the tunnel geometry as input and producing the trajectory of the molecule and energetic profile as output. The CaverDock workflow consists of multiple steps:

- 1) gather the input data (ligand in pdb or mol2 format, protein in pdb format), which can be obtained from experiments, downloaded from PDB² etc.
- 2) convert the input data into PDBQT format using AutoDock Tools [5]
- 3) identification and selection of a tunnel within the protein using Caver [8]
- discretization of the tunnel exported from Caver by the CaverDock script discretizer.py
- 5) (optional) setting the flexibility of selected sidechain residues by AutoDock Tools
- computing a box around the tunnel and the flexible residues either manually or using the CaverDock script prepareconf.py
- execution of CaverDock to search for the ligand trajectory
- analyze CaverDock trajectory and the energetic profile, and optionally identify new side chains which should be flexible and return to step 5

CaverDock's script flexibilize.py allows to automatically search for flexible residues. The script first runs CaverDock with the rigid receptor, and then iteratively runs CaverDock allowing the flexibility on the side-chains which have formed the bottlenecks in the previous iteration.

4 **TUNNEL DISCRETIZATION**

In this section we describe our requirements on the tunnel discretization in detail and the important parts of the algorithm performing the discretization.

2. https://www.wwpdb.org/

FILIPOVIČ, VÁVRA et al.: CAVERDOCK: A NOVEL METHOD FOR THE FAST ANALYSIS OF LIGAND TRANSPORT

4.1 Tunnel Discretization Requirements

As the first step in the CaverDock workflow, the tunnel must be discretized in discs, which will restrict the ligand's position in every docking. We use a geometric representation of the tunnel from Caver [8]. It approximates the tunnel as a sequence of spheres $T = {S_i}_{i=1}^k$, where the following holds:

$$1 \le i < k : S_i \bigcap S_{i+1} \ne \emptyset$$

$$\forall i \ne j : S_i \not\subseteq S_j$$
(3)

Moreover, the tunnel T never intersects itself, i.e. it is topologically equivalent to a cylinder.

Recall that the movement of the ligand is determined by the placement of its atom (the drag atom a_c) to the discs. Thus, we need to transform the tunnel T to a sequence of discs.

Definition 4.1. The cut θ of tunnel T is a disc in the threedimmensional space, which is defined by a triple $\theta = (A, u, r)$, where $A \in \mathbf{R}^3$ is a centre, $u \in \mathbf{R}^3$ is a normal and r > 0 is a radius. The $T \cap \theta$ must be a continuous set and $\exists \delta > 0$ such that $\forall \varepsilon > 0, \varepsilon < \delta$ holds $(A, u, r + \varepsilon) \cap T = \theta \cap T$.

Informally, Def. 4.1 ensures that a disc θ cuts the tunnel T in one place only, and it cuts it completely.

Having the discs cutting the tunnel defined, we can define how to generate cuts for the whole tunnel. Let $\Theta = \{\theta_i\}_{i=0}^n$ be a sequence of discs cutting tunnel T. We require the cuts to not intersect each other in more than a single point, formally:

$$\theta_i, \theta_j \in \Theta \Rightarrow |\theta_i \cap \theta_j| \le 1 \tag{4}$$

Moreover, we need to upper-bound the distance between discs, so we can also upper-bound the movement of the ligand atoms (to allow a contiguous trajectory generation). Let δ be an upper-bound of discs' distance and $\theta_i, \theta_{i+1} \in \Theta$. Then formally we require:

The fundamental requirement is to move forward in the tunnel, i. e. to generate a new cut ahead of the last cut:

$$\theta_i, \theta_{i+1} \in \Theta \Rightarrow \langle \theta_i^{normal}, \theta_{i+1}^{center} - \theta_i^{center} \rangle > 0$$
(6)

Finally, we want to start at the first sphere and end at the last sphere:

$$S_1^{center} \in \theta_1 \qquad \qquad S_k^{center} \in \theta_n$$
 (7)

4.2 Tunnel Discretization Computation

The discretization algorithm iteratively adds new discs to Θ . The tunnel geometry may be very complicated since the consecutive spheres may differ in radius significantly and may form sharp turns. Thus, we haven't found any simple analytical solution for the tunnel discretization. Instead, we have developed an iterative algorithm, which adds disc $\theta \in \Theta$ with the direction defined by a smoothed curve representing a tunnel and iteratively improve the positions of the disc to fulfill the requirements described in the previous section. In this section, we will focus on the main aspects of the algorithm, omitting the implementation details.



Fig. 4. Schematic 2D view of a curve (red) representing the tunnel direction. Left: naive solution where the centers of the spheres are connected; right: the centers of the minimal cuts (green) are connected.

4.2.1 Direction in the Tunnel

First, we define a curve, which represents a direction in the tunnel:

$$\gamma(t) \colon [0, l] \to \mathbb{R}^3 \tag{8}$$

5

where l is the length of the tunnel.

The easiest way to construct $\gamma(t)$ is to connect the centers of the spheres in T. However, some spheres may have a small influence on the shape of the tunnel and create a curve which will not represent the tunnel direction well (see Figure 4 left). Thus, we compute minimal cuts in the center of each sphere representing the tunnel and connect the centers of those cuts (see Figure 4 right). The minimal cuts are searched by an iterative optimization that uses the algorithm [20] for computing the smallest enclosing circle.

However, such construction is still not perfect, as it is not smooth. Let $C_1 \ldots C_n$ be the centers of tunnel's minimal cuts. For any $t_0 \in [0, l]$, we can easily find C_i, C_{i+1} such that t_0 lies in between them. We define a vector field $\Omega(t): [0, l] \to \mathbb{R}^3$. For a point t_0 , we define:

$$\Omega(t_0) = (1 - \lambda) \operatorname{norm}(C_{i+1} - C_i) + \lambda \operatorname{norm}(C_{i+2} - C_{i+1})$$
(9)

where

$$\lambda = \frac{\|\gamma(t_0) - C_i\|}{\|C_{i+1} - C_i\|} \tag{10}$$

if $i + 2 \le n$, else

$$\Omega(t_0) = \operatorname{norm}(C_{i+1} - C_i).$$
(11)

It represents a simple weighted average of vectors connecting centers of the minimal cuts. We further smooth $\Omega(t)$, creating a new vector field $\Phi : [0, l] \to \mathbb{R}^3$:

$$\Phi(t_0) = \int_{t_1 = \max\{t_0 - \Delta, 0\}}^{t_2 = \min\{t_0 + \Delta, l\}} \Omega(t) (\Delta - |t_0 - t|)^2 dt.$$
(12)

The vector field $\Phi(t)$ forms a smooth curve, which represents the direction in the tunnel well enough and is therefore used for initial placement of discs in the tunnel.

4.2.2 Helper Functions

Before we start with the tunnel discretization algorithm, we will introduce several important helper functions, which are used by the algorithm to place the discs along the curve $\Phi(t)$ representing the direction of the tunnel.

The function fitDiscTunnel computes the center and radius of the disc for the given plane ρ defined by the normal n and reference point P. The disc must be created to fulfill Definition 4.1: it must cut the tunnel at one place only and it must cut it completely. The function recursively builds a set of spheres $C \subseteq T$, which contain P or intersect

^{1545-5963 (}c) 2018 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information

6

IEEE/ACM TRANSACTIONS ON COMPUTATIONAL BIOLOGY AND BIOINFORMATICS

Algorithm 1 Algorithm for tunnel discretization

function DISCRETIZETUNNEL(T, δ) $centers \leftarrow [S^{center} \mid S \in \mathbf{T}]$ $discs \leftarrow [fitDiscTunnel(norm(S_1^{center} - S_0^{center}), S_0^{center})]$ $curve \leftarrow \text{TunnelCurve}(centers)$ $\begin{array}{l} \text{for } S_i \leftarrow S_0, \ldots, S_{|T|-2} \text{ do} \\ dir \leftarrow S_{i+1}^{center} - S_i^{center} \\ line \leftarrow \{S_i^{center} + t \cdot dir \mid t \in \mathbb{R}\} \end{array}$ $d \leftarrow 0$ while True do $prev_disc \leftarrow discs[|discs|-1]$ $plane \leftarrow \text{getPlane}(prev_disc)$ Construction of a plane containing disc $d \leftarrow \text{distance}(plane \cap line, S_i^{center}) + \epsilon$ if $d > \|dir\|$ then break end if end if if makesSharpTurn($prev_disc, curve$) then $disc_center \leftarrow prev_disc^{center} + \Delta * prev_disc^{normal}$ $disc_normal \leftarrow prev_disc^{normal}$ $doShift \leftarrow shiftSharpTurn$ else $\begin{array}{l} disc_center \leftarrow prev_disc^{center} + \epsilon * prev_disc^{normal} \\ disc_normal \leftarrow get Weighted Dir(curve, i, d) \end{array}$ $\mathrm{doShift} \gets \mathrm{shiftDisc}$ end if $disc \leftarrow fitDiscTunnel(disc_normal, disc_center)$ $disc \leftarrow doShift(prev_disc, disc)$ if $|discs| \ge 2 \land dst(disc, discs[|discs| - 2]) < \delta$ then $\operatorname{Pop}(discs)$ end if Append(discs, disc)end while end for return discs end function

both ρ and some sphere in *C*. Having the *C* constructed, the algorithm projects spheres from *C* to ρ and computes the circle encapsulating all projected spheres using the algorithm [20]. The computed circle determines the center and the radius of the computed disc, the normal of the disc is the same as the normal of ρ .

The function shiftDisc modifies the disc θ_{i+1} to fulfill Conditions 4, 5 and 6 in relation to the already placed disc θ_i . Let ρ be a plane orthogonal to planes where θ_i and θ_{i+1} lies. In ρ , discs θ_i , θ_{i+1} are projected as line segments. The algorithm is perfomed iterativelly: the line segment representing the disc θ_{i+1} is modified to not exceed δ in their ending points (Condition 5) and to not intersect (Condition 4). After that, disc θ_{i+1} is reconstructed from the projection and the fitDiscTunnel is called (as Definition 4.1 may be broken by shifting). This process is repeated untill there is no change on the disc θ_{i+1} .

In the case of a sharp curve in the tunnel, we need a more progressive placement of the disc, implemented in the function shiftSharpTurn. The plane θ_{i+1} is initially placed in Δ distance from θ_i , which breaks condition 5, but gives the idea of tunnel curvature. The function shifts the θ_{i+1} to intersect θ_i in the point nearest to θ_i , sets its center and normal to not exceed δ distance from θ_i and calls shiftDisc to finalize the θ_{i+1} placement. The main difference between shiftSharpTurn and shiftDisc is that shiftSharpTurn sets the initial position of θ_{i+1} such that its normal is pointing to the more distant direction of the tunnel.



Fig. 5. Discretization of a tunnel in the native toluene/o-xylene monooxygenase hydroxylase. The red circles represent the discs, the red arrows represent the tunnel direction and the grey balls represent the tunnel obtained from Caver [8].

4.2.3 Discretization Algorithm

The Algorithm 1 presents the main structure of the tunnel discretization algorithm. The input of the algorithm is a sequence of spheres T and the maximal distance between two discs δ . It uses helper functions, which are described in the previous section. The lines from 2 to 4 initialize the required data structures. The array *discs* will be used to construct Θ . During the initialization, the first disc is created with the same center as the first sphere and the normal given by the vector going from the center of the first to the center of the second sphere (so the first part of condition 7 is fulfilled). The *curve* contains the smoothed curve $\Phi(t)$ approximating the direction of the tunnel (see Section 4.2.1).

At line 5, the algorithm iterates over the spheres from T, where in each step it constructs *line*, containing the line connecting the actual and the next sphere. In the inner loop, the discs are generated from the centre of the sphere S_i to S_{i+1} (condition at line 13). The variable *d* determinates our distance from S_i^{center} and controls the number of loop iterations at line 9.

The function isSharpTurn detects if the tunnel forms a sharp turn at the particular place. If so, the algorithm uses a more aggressive strategy for the next disc placement: it places a new disc parallel to the current disc with distance Δ . Otherwise, the center of the new disc is displaced by ϵ , and its normal is set according to the weighted direction curve. The constant Δ determines the distance, which is checked for deciding whether the turn is sharp. We use $\Delta = 2\delta$. The constant ϵ determines the granularity of discretization, and it must be lower than δ . We set $\epsilon = \frac{1}{10}\delta$.

After the initial disc placement the function fitDiscTunnel is called. The function computes the center and radius of the disc, so the disc forms the cut of the tunnel (see Definition 4.1), and the radius of the disc is minimal. After fitting the disc, the function doShift may further improve its placement according to the curvature of the tunnel, so the disc will be placed to fulfill Conditions 4, 5 and 6. Finally, the algorithm checks if the disc created in the previous iteration can be omitted (to not generate too dense discretization) at lines from 27 to 28. An example of the algorithm output is depicted in Figure 5.

5 DOCKING WITH RESTRAINTS

As we have described in Section 3.2, we are employing two types of restraints. Recall that the tunnel-position restraint

FILIPOVIČ, VÁVRA et al.: CAVERDOCK: A NOVEL METHOD FOR THE FAST ANALYSIS OF LIGAND TRANSPORT

snaps a selected atom a_c in ligand λ to the disc θ and the pattern restraint places the ligand λ in the vicinity of $\lambda_{pattern}$.

During the docking, the position of a protein-ligand complex is minimized. Therefore, we need to impose spatial restraints in the form of energy terms that are added to the force field function from AutoDock, that will be minimized when searching for the optimal docking position within each disk. However, the penalty produced by restraints is used to hold a ligand in a defined area during minimization and is not added to the AutoDock Vina energy after minimization. Therefore, only the original AutoDock Vina energies are reported as the binding free energy of the protein-ligand complex.

In this section, we first introduce the search-space optimization methods implemented in AutoDock Vina and after that describe how the newly-added restraints are implemented.

5.1 AutoDock Vina Search Space Methods

The molecular docking is an optimization problem, where the docking program is searching for the global minimum of energy defined by the position of the ligand and the flexible parts of the receptor with respect to the given force field. Two optimization methods are working together in AutoDock Vina: a stochastic global optimization and a gradient-based local optimization.

The simulated system has multiple degrees of freedom (DoF), which must be searched. First, the ligand is considered as a body in the space, having its position and orientation vectors (six dimensions). Second, the ligand is a flexible body – it may be bent by setting angular values for its free dihedral angles (one dimension per dihedral angle on every single bond). Third, the receptor may contain flexible side-chains (so also the receptor geometry may be partially flexible), where each flexible residue contains one or more free dihedral angles. Thus, the optimization algorithm must optimize a high number of DoF (typically tens). The nature of the chemical force field creates a lot of local minima in the search space.

The global-optimization method implemented in AutoDock Vina is based on the Markov chain Monte Carlo method (MCMC). The initial state (position and orientation of the ligand, dihedral angles of the ligand and flexible side chains) is selected randomly keeping the ligand within the defined box. Subsequently, a predefined number of globaloptimization steps are performed. In a global optimization step, one or more DoF are changed by an upper-bound random value, so the newly-generated conformation of the ligand and flexible side-chains is in the upper-bound vicinity of the previous conformation. After the global optimization step, the local optimization is executed. If the local optimization converges to a better minimum than what was reachable from the previous global-optimization step, the global-optimization accepts the new step and uses it as a base for the next iteration. Otherwise, the new step is accepted only with small probability based on the Metropolis critorion (this feature allows the algorithm to escape from a local minimum). During the global-optimization, the significant local minima are stored, so AutoDock Vina is able to return multiple different conformations, not only the best one.

The local-optimization in AutoDock Vina implements the gradient-based Broyden-Fletcher-Goldfarb-Shanno (BFGS) method [21]. It is a variant of the Newton method, so derivatives of all force-field terms have to be computed (it also uses the second derivatives, but they are computed numerically).

5.2 Tunnel-position Restraint

To snap atom $a_c \in \lambda$ to the disc θ , we add a force-field term which penalizes the ligand's positions where $|a_c - t| > 0$, where $t \in \theta : \forall u \in \theta, u \neq t, |a_c - u| > |a_c - t|$ (i.e. the distance is computed as the distance between a_c and the nearest point in θ). The penalization energy e_p is computed as a Gauss function of the distance $|a_c - t|$:

$$p_{position} = p_{max} - p_{max} e^{-\frac{|a_c - t|^2}{0.5}}$$
 (13)

where p_{max} is the maximal value of the penalization energy. The constant 0.5 used in the exponent has been selected experimentally. It ensures that the half of e_{max} penalization is applied when $|a_c, t| = 0.5$ Å. Note that the bell-shaped function is used to avoid strong penalization of small distances between a_c and θ too strongly in order to to keep the good numerical stability of the BFGS optimization method (so the Condition 1 can be violated by a small distance in practice).

The term in Eq. 13 and its derivative has been added into the energy and force evaluation codes in AutoDock Vina, so it is applied to the dragged atom a_c during the BFGS local optimization. Moreover, we have added a simple modification into the MCMC global optimization: when a new conformation is randomly generated, the ligand is shifted by a vector $t - a_c$, so the global optimization step does not break the tunnel-position restraint. Note that the modification of the global optimization method is not necessary for applying the restraint in the docking, however, it speeds up the docking convergence.

5.3 Pattern Restraint

The pattern restraint keeps the ligand λ in the vicinity of the pattern position $\lambda_{pattern}$ (Equation 2), so it must be applied to all atoms of the ligand. The pattern restraint is also applied to the flexible side chains, however, for the sake of simplicity, we describe the application to a ligand only (the principle of the pattern is the same for the ligand and flexible side chains).

The pattern restraint is applied for all pairs of corresponding atoms $a \in \lambda$ and $b \in \lambda_{pattern}$. Let δ be the distance which is not penalized by the pattern restraint. The energy of the pattern restraint is computed as:

$$e_{pattern} = c \cdot \sum_{a \in \lambda, b \in \lambda_{pattern}} \max\left(0, |a - b| - \delta\right)$$
(14)

where *c* is a constant determining the strength of the pattern (it has been empirically set to 40). Apparently, Eq. 14 is not differentiable in the area where $|a-b| = \delta$. We define a derivative at these points to be 0 and keep the computation

8

IEEE/ACM TRANSACTIONS ON COMPUTATIONAL BIOLOGY AND BIOINFORMATICS

of the pattern restraint simple for the sake of computational efficiency.

The MCMC global optimization method is constructed to perform a long chain of conformational changes to escape from the local minima. However, when the pattern restraint is applied, the movement of the ligand is restricted to the vicinity of the pattern. Thus, we have modified the global search, such that (i) the initial configuration mimics the position of the pattern and (ii) the number of steps of the global optimization is $100 \times$ lower compared to the default setup. The MCMC method still allows to escape from the local minima but does not generate too long chains due to the limited range of ligand movements from the initial configuration.

Note that, as all restraints are evaluated as force field terms, they can be violated if there is some strong energy contribution generated by different force field term (e.g. the pattern restraint may be violated when pushing the ligand against a rigid part of the receptor). Therefore, CaverDock filters the computed conformations and discard those that where the restraint violations exceeding some threshold (e.g. if we consider contiguous conformation changes up to 0.5 Å, we can set the pattern restraint to penalize movement larger than 0.4 Å and tolerate the movement not exceeding 0.5 Å).

6 TRAJECTORY SEARCH

The implemented restraints allow us to define the ligand's position in the tunnel and upper-bound its distance from some pattern. Thus, it is possible to iteratively dock the ligand along the tunnel and analyze the energy of the transport process. Recall that we compute two types of trajectory:

- lower-bound trajectory, which samples the tunnel finely, but the movement of the ligand and flexible side-chains is not contiguous;
- upper-bound trajectory, which is contiguous.

The lower-bound trajectory may underestimate the energy of barriers, as the ligand may flip or change its conformation dramatically between two consecutive steps (Figure 1). However, its computation is straightforward – there is no dependence between consecutive steps (only the tunnel-position restraint is used), and therefore we may perform only n docking steps, where n is the number of discs. The upper-bound trajectory generates a contiguous movement of the ligand and side-chains using the pattern restraint. However, a high number of possible contiguous trajectories exist and there is no guarantee that our method finds the lowest energy trajectory. Thus, we call the trajectory upperbound, as it is not known if its energy may be further improved. The optimal energies should lie between upper-and lower-bound values.

Having a set of discs $\theta_1 \dots \theta_n$, the lower-bound trajectory is defined as

$$\Lambda_{lb} = \lambda_{min}^1, \lambda_{min}^2, \dots, \lambda_{min}^n \tag{15}$$

where λ_{min}^{i} denotes the conformation at disc *i* with the lowest energy from all explored λ^{i} .

$$\Lambda_{ub} = \lambda_1^1, \dots, \lambda_{m_1}^1, \lambda_1^2, \dots, \lambda_{m_2}^2, \dots, \lambda_1^n, \dots, \lambda_{m_n}^n$$
 (16)

where $m_1 \dots m_n \ge 1$, so the upper-bound trajectory follows a forward movement within the tunnel or changes the ligand position on a disc, but does not go backward.

6.1 Ligand Movement Driving

The trajectory search is driven by a set of final state automata. Each automaton is designed to perform different tasks:

- general automaton, controlling the overall progress of the trajectory search;
- *lower-bound trajectory automaton,* performing lowerbound trajectory computation;
- *forward movement automaton,* responsible for moving the ligand forward in the tunnel;
- *optimization automaton,* optimizing the position of the ligand at the particular disc;
- *backtracking automaton*, moving the ligand backward if it hits a barrier;
- convergence automaton, pushing the ligand in a backtracked trajectory to converge with the forward trajectory

The execution of automata can be nested. For example, the general automaton calls the forward automaton to move forward in the tunnel, and the forward automaton calls the optimization automaton to improve the position on a disc. The reason for such an implementation is twofold. First, the different operations on the trajectory are separated and the code is easier to maintain. Second, the computation can be interrupted or altered at any place since each automaton performs at most one state transition per call. Thus, it is, for example, possible to execute multiple backtracking in parallel, as the general automaton can immediately continue after spanning a new backtracking automaton without waiting for the backtracking to finish.

The automata call all the restrained docking computations in a non-blocking manner. They submit tasks into an internal CaverDock queue, which is then processed in parallel. Therefore, it is possible to process multiple alternative trajectories in parallel by executing multiple automata in a simple serial loop, or an automaton may process multiple alternatives at once.

6.1.1 Trajectory Search Parameters

The algorithm for trajectory search uses several parameters, affecting its precision and a number of executed dockings (and hence the computation time). Those parameters, listed below, affect the state machines or the docking settings and may be configured by the user.

- Parameter *optimization stratregy* determines the optimization criterion for the trajectory search. In the current implementation, we can execute CaverDock to minimize the highest energy peak across the whole trajectory or minimize the integral of the trajectory energy.
- Parameter backtrack threshold quantifies the energy difference (in kcal/mol) between the lower-bound

^{1545-5963 (}c) 2018 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

FILIPOVIČ, VÁVRA et al.: CAVERDOCK: A NOVEL METHOD FOR THE FAST ANALYSIS OF LIGAND TRANSPORT



Fig. 6. General automaton.

and upper-bound energy which triggers the backtracking. It is expectable that the energy of contiguous upper-bound trajectory will be higher, however, too large difference may indicate that the upperbound trajectory is suboptimal. Therefore, when the difference between the upper-bound and lowerbound trajectory energies at some disc exceeds the threshold, the backtracking is used to search for a better trajectory.

- Parameter *backtrack limit* sets the number of discs that are processed before a new backtracking can be executed. The parameter may speed-up CaverDock when the number of executed backtrackings is too high. This parameter is ignored when the forward trajectory cannot be computed because of a bottleneck (the backtracking is then started immediately).
- Parameter *contiguous threshold* sets the highest distance which the atoms can move between consecutive conformations, if this movement is considered contiguous.
- Parameter *pattern limit* sets the highest distance that is not penalized by the pattern restraint. The pattern limit must be lower than the contiguous threshold so the pattern restraint may actually apply some force to ligand position before the ligand position is discarted.

6.1.2 General Automaton

The simplified scheme of the general automaton is shown in Figure 6. After initialization, it starts to compute a lower-bound trajectory Λ_{lb} and builds a cache of alternative conformations Λ_{cache} (all examined conformations on discs $1 \dots n$). When the lower-bound computation is not successful (i. e. $\exists i \in <1, n >, \lambda^i \notin \Lambda_{lb}$), the automaton halts in a *LB failed* state. It may happen when the tunnel is very narrow in some part, and it is not possible to dock the ligand there. Otherwise, the general automaton starts with searching for a contiguous upper-bound trajectory. It inserts λ_{min}^1 into Λ_{ub} and moves it into the *forward* state performing the following steps:

• Call the forward automaton to move forward from the last position $\lambda^i \in \Lambda_{ub}$ till it reaches the end of the tunnel, or requests backtracking. If the forward automaton requests backtracking (the energy of the forward trajectory is too high comparing to the lower-bound), create a backtracking automaton and change the state to backtrack. The backtracking automaton starts from conformation $\lambda_i^{i+1} \in \Lambda_{cache}$, where *i* is the last position in the forward trajectory Λ_{lb} and j is selected such that λ_i^{i+1} has not been used for backtracking so far and its energy is minimal. It builds a backtrack trajectory $\Lambda_{backtrack}$. If the trajectory $\Lambda_{backtrack}$ is successfully found and improves the energy of the trajectory, it is implemented into Λ_{ub} . More precisely, the conformations in Λ_{ub} , from the conformation where the backtracking trajectory can be connected to the end of the trajectory, are removed and then $\Lambda_{ub} \leftarrow \Lambda_{ub} \cup \Lambda_{backtrack}$. Otherwise, a new backtracking is executed using a different starting conformation from Λ_{cache} , which has not been used for backtracking so far. If no such a conformation exists, then the general automaton returns to the forward state.

9

• If the forward automaton requests a forced backtracking (it cannot find a forward trajectory), the backtracking automaton is created as in the previous case, and the general automaton changes its state to *forced backtrack*. The difference to the *forced backtrack* state is that $\Lambda_{backtrack}$ is implemented into Λ_{ub} every time when it is found. If $\Lambda_{backtrack}$ cannot be found, the general automaton ends in *UB failed* state: the upper-bound trajectory cannot be computed.

6.1.3 Other Automatons

The forward and backward automatons are responsible for moving the ligand in the tunnel. They perform essentially two steps: transition to a different disc and optimization of ligand's position on the same disc. The optimization is performed by the optimization automaton, which searches for trajectory $\lambda_2^i \in \Delta \lambda_1^i, \lambda_3^i \in \Delta \lambda_2^i \dots$ until the energy of ligand is improved. With backtracking automaton, optimization is replaced by convergence every five steps. The convergence automaton optimizes the ligand position at the same disc similarly to the optimization automaton. However, instead of moving the ligand to the local minimum, it uses a soft pattern restraint to attract the ligand atoms to a position determined by $\lambda_{dest}^i \in \Lambda_{ub}$ (by setting the constant c in Equation 14 to one, instead of 40 used in the restraint forcing the contiguous movement). Thus, it forces backtracking trajectory to converge with the forward one.

6.2 Software Architecture

The CaverDock is built as an MPI application using masterslave parallelism. There is one master process, driving the trajectory search (i. e. executing automatons and assigning work for slaves). The slave processes are responsible for computing the restrained docking: they receive restraints from the master (position of the disc and position of pattern atoms) and send the computed conformations with the computed energies (chemical force field and restraints' energy).

The master process runs in a loop, querying automatons and gathering data from slaves. The automatons are called

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCBB.2019.2907492, IEEE/ACM Transactions on Computational Biology and Bioinformatics

IEEE/ACM TRANSACTIONS ON COMPUTATIONAL BIOLOGY AND BIOINFORMATICS

my

Fig. 7. Arachidonic acid.

10

in a non-blocking fashion: they submit a work package describing the input for the docking. This work-package is held by the master process and assigned to a slave when is ready. Therefore, automatons can submit any number of work packages without waiting for the result and react by changing its state when the work is completely done.

The CaverDock can also be executed in a simple docking mode. In such case, only one slave process is executed, taking the input for the docking from a command line. When the user executes CaverDock with the parameters used in the original AutoDock Vina, CaverDock operates exactly as AutoDock Vina. However, it is possible to also pass the restraint via command line. Therefore, CaverDock is usable also as a docking tool allowing richer control over the docking process via restraints. It may be applied for observing a particular docking conformation when the user is interested in searching for a ligand conformation with some atoms restricted in a defined area of their interest.

7 EVALUATION

In this section we compare the results of CaverDock with similar tools and demonstrate CaverDock's ability to analyze complex tunnels in reasonable time on chemicallyrelevant data. The limited evaluation showing chemical relevance of the computed results is demonstrated on two cases: ligand unbinding and reproduction of ligand positions. The comprehensive evaluation of CaverDock is being presented in parallel in other papers [6], [7].

7.1 Testbed Setup

For testing the stability and time demands of CaverDock, we have prepared a representative set of biologically relevant protein-ligand pairs shown in Table 1. The set contains proteins with both short and long tunnels (e.g. the insulin hexamer tunnel is discretized to 42 discs only, whereas the glucose transporter tunnel is discretized to 362 discs). The complexity of ligands also heavily varies: phenol has only 7 DoF (6 for the position and orientation and 1 free dihedral angle) whereas the arachidonic acid has 20 DoF (14 free dihedral angles).

We have tested CaverDock runtime using desktop computer equipped by AMD Ryzen 7 1700 (8 cores at 3.0 GHz) and 16 GB RAM. The resulting time, quality of the result and a number of performed docking calculations are shown in Table 2.

7.2 CaverDock Runtime and Robustness

CaverDock was not able to compute the upper-bound trajectory in two cases (vitamin D receptor + 1,25-dihydroxyvitamin D3 and cytochrome P450 2E1 + arachidonic acid), whereas the lower-bound trajectory has been

TABLE 1 The experimental set of molecules used for CaverDock evaluation.

protein+ligand	ligand DoF	discs
haloalkane dehalogenase + 1-chlorobutane	8	72
acetylcholinesterase + acetylcholine	10	85
leucine transporter + leucine	10	105
lactose permease + lactose	18	138
glucose transporter + glucose lipase B + 4-methyloctanoic acid insulin hexamer + phenol	12	362
	12	60
	7	42
aquaporin Z + glycerol	11	121
vitamin D receptor + 1,25-dihydroxyvitamin D3 cytochrome P450 2E1 + arachidonic acid	13	92
	20	149

computed in all of the tested cases. The arachidonic acid contains a long chain with a high number of DoF (see Figure 7), which complicates the process of searching for the upper-bound trajectory. We suppose that CaverDock heuristics fails to find the contiguous movement of such a complicated molecule. 1,25-dihydroxyvitamin D3 has also a high number of DoF, but the main reason why CaverDock failed to compute contiguous trajectory is due to the narrow part of the tunnel entrance, which is difficult to pass with a contiguous movement. In contrast, lactose has also a high number of DoF, but the tunnel in the lactose permease is wider and CaverDock had no problem to compute the contiguous upper-bound trajectory.

The computation time ranges from 1 m3 s to nearly 2.5 h. The CaverDock heuristic is in $O(n^2)$, where *n* is the number of discs (as backtracking can be issued during the whole trajectory and, in the worst case, may continue to the trajectory beginning). Therefore, the number of docking calculations, and hence the computational time, may grow quadratically with the tunnel length. However, the backtracking is not issued often in narrow tunnels, where the running time may grow according to the number of discs. The time required for each docking is highly dependent on the number of DoF, for example, 14.3 dockings per second are computed in the case of phenol (7 DoF), but only 1.33 dockings per second are computed in the case of lactose (18 DoF).

7.3 Energy Profiles

The energy profiles of the tested ligand-protein complexes are shown in Figure 8. They represent the variation of the binding energy of the ligands moved from the active site to the tunnel entrance at the protein surface.

In all experiments the upper-bound trajectory has higher energy than the lower-bound because the ligand movement is restricted by the pattern restraint, and thus cannot easily overcome small bottlenecks. Whereas the upper-bound

FILIPOVIČ, VÁVRA et al.: CAVERDOCK: A NOVEL METHOD FOR THE FAST ANALYSIS OF LIGAND TRANSPORT

TABLE 2
The output characteristics and computational demand of CaverDock
calculations using ten different biological systems. LB = lower-bound,
UB = upper-bound.

protein	result	runtime	dockings
haloalkane dehalogenase	LB+UB	2 m 14 s	1,824
acetylcholinesterase	LB+UB	4 m 46 s	1,920
leucine transporter	LB+UB	7 m 4 s	2,688
lactose permease	LB+UB	67 m 20 s	5,384
glucose transporter	LB+UB	149 m 47 s	23,888
lipase B	LB+UB	3 m 19 s	896
insulin hexamer	LB+UB	1 m 3 s	900
aquaporin Z	LB+UB	8 m 50 s	3,896
vitamin D receptor	LB only	40 m 6 s	4,132
cytochrome P450 2E1	LB only	19 m 35 s	644

energies usually copy the shape of the lower-bound energetic profile, some bottlenecks are visible in upper-bound trajectory only (such as in acetylcholinesterase from distance 7 Å to 16 Å, in lipase B from 3 Å to 6 Å and aquaporin Z from 13 Å to 18 Å). As we can see, the contiguous trajectory adds additional information useful for the ligand transport analysis. Note that the observed bottleneck does not necessarily indicate that the transportation of a ligand through a protein tunnel is not possible. The protein flexibility may allow the ligand to pass even higher energy barriers observed in static structures. The interpretation of such data is crucial - the part of the protein forming the bottleneck may be more or less rigid in a real-world system, or sometimes some form of flexibility may lead to the opening of the tunnel and allow the ligand to pass. The CaverDock user may select the residues to which the side chain flexibility may be introduced. If a protein backbone forms an artificial bottleneck, then a different protein conformation has to be used.

7.4 Comparison with Similar Tools

We have tested the set of molecules described in Table 1 also with SLITHER [19] and MoMA-LigPath [13]. We were not able to compute the trajectories in the tunnels of lipase B, insulin hexamer, aquaporin Z, vitamin D receptor and cytochrome P450 2E1 with SLITHER, and in glucose transporter, aquaporin Z, vitamin D receptor and cytochrome P450 2E1 with MoMA-LigPath. Therefore, at least for our testing set, CaverDock was more robust regarding its ability to compute the trajectories. The runtime of SLITHER and MoMA-LigPath is in order of minutes in the worst case. Therefore, CaverDock time is comparable for simpler cases but may be significantly higher when a large number of dockings needs to be executed.

The trajectories produced by CaverDock qualitatively differs from the trajectories obtained with SLITHER and MoMA-LigPath. More precisely, SLITHER generates scattered, non-contiguous trajectory, whereas MoMA-LigPath does not produce energy profiles. We demonstrate the difference in the produced trajectories using an example of the transportation of acetylcholine through a tunnel in the protein acetylcholinesterase (PDB ID 1MAH). The acetylcholine has been moved through the tunnel from the active site to the protein surface. The trajectory computed by CaverDock



11

Fig. 8. Energy profiles of the ligand movements in the test set defined in Table 1. The distance is measured from the tunnel bottom to the protein surface.

is shown at Figure 9. It can be seen that there are no gaps (empty spaces) in acetylcholine trajectory – the movement of its atoms is contiguous.

The trajectories computed by SLITHER and MoMA-LigPath are shown in Figure 10. SLITHER does not implement a restrained docking, and as a consequence it produces large gaps in the computed trajectory. MoMA-LigPath, on the other hand, produces a contiguous trajectory. However, no chemical force field is used in MoMA-LigPath, so the user has no information describing the energy profile associated with the trajectory. Moreover, the trajectory does not reflect the chemical interactions and therefore can follow a path which would not be favored in the real systems. It can also be seen that the trajectory produced by MoMA-LigPath is more regular compared to CaverDock (the direction of atoms' movements is similar in multiple following conformations), which is very likely a consequence of the missing chemical forces, which increase the complexity of

12



Fig. 9. Trajectory of acetylcholine in the tunnel of acetylcholinesterase computed by CaverDock. All 174 positions of acetylcholine are shown as the superimposed cyan sticks. The surface of the protein (close-up of tunnel) is shown as the grey surface, the catalytic S203 ($C\alpha$ atom) as the magenta sphere, and the chemical structure of this ligand is shown on the right.



Fig. 10. Trajectory of acetylcholine in the tunnel of the acetylcholinesterase computed by SLITHER (left image) and MoMA-LigPath (right image). The surface of the protein (close-up of tunnel) is shown as the grey surface, the catalytic S203 (C atom) as the magenta sphere.

the optimized trajectory but are essential to describe the realistic behavior of the molecules.

Note that although visualized for acetylcholinesterase (Figure 9 and 10), SLITHER has generated trajectories with similar gaps also for other test cases. We have compared the CaverDock energy profiles to SLITHER in cases where we were able to compute the SLITHER trajectory (Figure 8). The data provided by SLITHER has been filtered: we have removed all conformations which were not placed within the tunnel determined by CAVER. Such filtering is necessary to exclude any conformations in different tunnels or at the protein surface. The advantage of the restrained docking used in CaverDock can be seen when the energy profiles are compared. The trajectory obtained with SLITHER was sparser when compared to CaverDock and no conformation was placed in the bottleneck. For example, only two conformations were computed in the leucine transporter's tunnel and they are located before and after the bottleneck. In some cases it might be possible to roughly guess the position of the bottlenecks based on the gaps in SLITHER's trajectory. In other cases, SLIGHTER's trajectory may include gaps also at the positions with no observable bottleneck, which can be seen for example in the second half of the trajectory with haloalkane dehalogenase. CaverDock reports the energies and conformations of the ligand in the bottlenecks, so it is possible to analyze, which residues may be mutated to increase the rate of ligand's passage. Note that the absolute energy values are different for SLITHER and CaverDock, which is caused by using different chemical force fields (the restraints force field terms are excluded from CaverDock output energies).

IEEE/ACM TRANSACTIONS ON COMPUTATIONAL BIOLOGY AND BIOINFORMATICS

7.5 Protein-ligand Unbinding Test Case

CaverDock was tested in a biotechnologically relevant case study, published recently [22]. It is known that the release of the product 2,3-dichloropropan-1-ol (DCP) from the buried active site is the rate-limiting step in the catalytic conversion of 1,2,3-trichloropropane (TCP) into DCP by the haloalkane dehalogenase DhaA31. Free energy calculations were performed after exhaustive metadynamics simulations to determine the free energy profiles of DCP traveling through the tunnels of DhaA31 and the wildtype DhaAwt. It was found that the energy barrier to the release of DCP was higher in DhaA31 than in DhaAwt by 2.5 kcal/mol. This explained why DhaA31 is less prone to release DCP than DhaAwt, which was in agreement with all the evidence. When CaverDock was used to calculate the energy profiles of DCP through the tunnels of those enzyme variants, the energy barrier was much higher for DhaA31 than for DhaAwt, which followed the same trend as the free energy results [22]. The energy barriers were always located at the tunnel bottlenecks. The binding energy values calculated with CaverDock were, however, different from the respective free energy profiles. This was expected for several reasons: 1) those CaverDock simulations did not take into account the flexibility of the receptor during the transport of the ligand, which necessarily raises the energy with unnatural clashes; 2) the binding energy calculated by Autodock Vina is only one component of the free energy, and does not take into account the entropy of the system; 3) the solvation of the ligand may influence the energy and facilitate the release. Nonetheless, CaverDock allowed to correctly predict which enzyme can release DCP faster, and helped in the identification of some of the residues in the tunnels that interacted more strongly with DCP and that may prevent its release. This information can be useful for the design of improved biocatalysts. The metadynamics simulations needed several weeks to be completed, while each CaverDock calculation was performed in less than 0.5 hours.

7.6 Reproduction of Ligand Positions Determined by Protein Crystalography

The ability of CaverDock to reproduce the protein-ligand complexes determined by protein crystallography was tested with the crystal structures of protein Hsp90 bound with 34 inhibitors [23]. The list of PDB IDs with the respective results is provided in the Table 3.

We calculated the root mean squared deviations (RMSDs) between the positions of the inhibitors found in the crystal structures with the snapshots obtained from Caver-Docks lower-bound trajectories, and the docked poses from AutoDock Vina (the exhaustiveness was set to 10). In Table 3 we report both the lowest RMSDs as well as the RMSDs for the lowest energy conformations from CaverDock and AutoDock Vina. In most cases we reached low RMSD values (≤ 2.0 Å) fitting well the experimental data, and were as precise or better than classical docking. CaverDock was able to find significantly better complex than the docking in 9 cases out of 34 (PDB ID 2VCI, 5J2X, 6F1N, 6ELN, 5J86, 5LQ9, 5ODX, 6EL5 and 5LO5). CaverDock was able to get spatially close to the position of the original inhibitors, but the lowest

FILIPOVIČ, VÁVRA et al.: CAVERDOCK: A NOVEL METHOD FOR THE FAST ANALYSIS OF LIGAND TRANSPORT

TABLE 3 RMSD (in Å) of complexes obtained by CaverDock and AutoDock Vina compared to complexes obtained by crystallography.

	The Closest Pose		The Lowest Energy Pose	
PDB ID	Docking	CaverDock	Docking	CaverDock
2VCI	9.20	0.86	10.09	1.25
2UWD	0.70	0.67	0.70	0.67
2BSM	0.68	0.60	0.68	0.60
5NYI	1.41	0.60	1.41	2.01
5J2X	5.90	0.83	6.84	1.98
6F1N	7.04	0.81	7.27	1.87
6ELO	0.37	0.31	0.37	0.31
5J64	1.53	0.54	1.56	1.56
6ELN	1.18	0.49	1.18	0.49
5J20	0.56	0.50	0.56	8.30
5J86	3.43	0.46	6.51	0.49
5J9X	0.76	0.66	0.76	0.67
6ELP	0.44	0.26	0.44	0.29
5J27	1.37	0.41	1.37	0.45
5LRZ	0.28	3.14	2.03	3.58
5LR7	0.45	9.83	0.45	9.88
2YKI	8.64	2.65	9.72	9.57
5LQ9	7.29	0.80	10.55	0.84
5LS1	1.14	0.84	1.14	1.12
5T21	1.21	0.94	1.21	1.31
6EYA	1.17	1.08	1.17	1.08
5LO6	0.92	0.61	0.92	0.61
5LNZ	0.91	1.10	0.91	1.31
6EY8	0.45	1.15	0.45	5.49
6EY9	0.95	0.39	0.95	1.29
50CI	1.58	1.26	1.58	1.51
50DX	0.77	1.07	6.24	1.10
5NYH	0.44	0.44	0.44	0.44
50D7	0.53	0.48	0.53	0.51
6EI5	0.65	0.44	0.65	0.77
5LR1	0.22	0.42	0.22	0.42
6EL5	3.66	0.48	3.70	1.71
5LO5	1.40	1.13	4.49	1.13
2YKJ	0.37	1.61	0.37	9.32
Avg.	1.99	1.11	2.57	2.17

energy conformation was located in a different part of the trajectory in 4 cases (PDB ID 5J20, 5LR7, 6EY8 and 2YKJ). CaverDock failed in finding the correct conformation for the closest and the lowest energy case for the complex 5LR7. The docking did not reproduce the correct binding pose in this case due to the large size of the search space, which might be improved by setting a higher exhaustiveness. The high RMSDs from CaverDock may be caused by incorrect orientation of the ligand and also by the location of the original inhibitor, which was far from the tunnel. The positions of the inhibitors buried deep in the protein structure, outside the access tunnels, may not be reachable by CaverDock since the ligand is always spatially constrained to the disks.

8 CONCLUSIONS AND FUTURE WORK

In this paper, we have introduced a novel method for analysis of the transport of ligands in proteins and its implementation in CaverDock tool. We have developed a restrained docking and a heuristics to analyze the ligand movements in the tunnel. We have also introduced a new algorithm for discretization of the protein tunnels. Our approach extends the state-of-the-art by using molecular docking for calculation of contiguous movements of the ligand within the tunnel. The calculation is faster and easier to setup compared to MD but on the other hand overcomes the limitations of geometrical methods. We have demonstrated that CaverDock is robust and is able to analyze the ligand transportation usually in minutes, or in a few hours in the worst scenario.

13

In the future, we plan to improve CaverDock heuristics to compute more alternative trajectories and to connect promising, non-contiguous parts of trajectories more aggressively. We expect to obtain lower energy for upper-bound trajectories. Or, at least, generate upper-bound trajectories with higher confidence, so the energy is not overestimated due to insufficient sampling. Furthermore, we plan to improve the receptor flexibility in CaverDock. With the current version, only the side-chains can be flexible. The flexibility of the protein backbone would allow to model situations where the receptor's flexibility plays a significant role in the ligand passage. More precisely, we will explore the possibility to use an ensemble of protein conformations or coarse-grained MD to reproduce the movement of the protein backbone.

ACKNOWLEDGMENTS

The work was supported from Grant Agency of Masaryk University (MUNI/M/1888/2014) and European Regional Development Fund, Project "CERIT Scientific Cloud" (No. CZ.02.1.01/0.0/0.0/16_013/0001802). The authors express their thanks also to infrastructural projects ELIXIR and C4Sys run by the Czech Ministry of Education (LM2015047 and LM2015055) for financial support. O. Vávra is the recipient of a Ph.D. Talent award provided by Brno City Municipality. Access to the CERIT-SC computing and storage facilities provided by the CERIT-SC Center, under the programme "Projects of Large Research, Development, and Innovations Infrastructures" (CERIT Scientific Cloud LM2015085), is greatly appreciated.

REFERENCES

- S. Marques, L. Daniel, T. Buryska, Z. Prokop, J. Brezovsky, and J. Damborsky, "Enzyme tunnels and gates as relevant targets in drug design," *Medicinal Research Reviews*, vol. 37, no. 5, pp. 1095– 1139, 2017.
- [2] O. Trott and A. J. Olson, "Autodock vina: Improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading," *Journal of Computational Chemistry*, vol. 31, no. 2, 2010.
- [3] T. J. A. Ewing, S. Makino, A. G. Skillman, and I. D. Kuntz, "DOCK 4.0: Search strategies for automated molecular docking of flexible molecule databases," *Journal of Computer-Aided Molecular Design*, vol. 15, no. 5, 2001.
- [4] R. Thomsen and M. H. Christensen, "MolDock: a new technique for high-accuracy molecular docking," *Journal of Medicinal Chemistry*, vol. 49, no. 11, pp. 3315–3321, 2006.
- [5] G. M. Morris, R. Huey, W. Lindstrom, M. F. Sanner, R. K. Belew, D. S. Goodsell, and A. J. Olson, "AutoDock4 and AutoDockTools4: Automated docking with selective receptor flexibility," *Journal of Computational Chemistry*, vol. 30, no. 16, 2009.
- [6] O. Vávra, J. Filipovič, J. Plhák, D. Bednář, S. M. Marques, A. Pavelka, L. Matyska, and J. Damborský, "CaverDock: Ligand transport analysis based on molecular docking," 2018, in preparation.
- [7] G. Pinto, O. Vávra, J. Filipovič, D. Bednář, and J. Damborský, "Fast screening of binding and unbinding of inhibitors using novel software tool CaverDock," 2018, in preparation.

14

IEEE/ACM TRANSACTIONS ON COMPUTATIONAL BIOLOGY AND BIOINFORMATICS

- [8] E. Chovancová, A. Pavelka, P. Beneš, O. Strnad, J. Brezovský, B. Kozlíková, and et al., "CAVER 3.0: A tool for the analysis of transport pathways in dynamic protein structures," *PLoS Computational Biology*, vol. 8, no. 10, 2012.
- [9] D. Sehnal, R. Svobodová Vařeková, K. Berka, L. Pravda, V. Navrátilová, P. Banáš, C. M. Ionescu, M. Otyepka, and J. Koča, "MOLE 2.0: advanced approach for analysis of biomacromolecular channels," *Journal of Cheminformatics*, vol. 39, no. 5, 2013.
- [10] E. Yaffe, D. Fishelovitch, H. J. Wolfson, D. Halperin, and R. Nussinov, "MolAxis: Efficient and accurate identification of channels in macromolecules," *Proteins: Structure, Function, and Bioinformatics*, vol. 73, no. 1, pp. 72–86, 2008.
 [11] J. Brezovsky, E. Chovancova, A. Gora, A. Pavelka, L. Biederman-
- [11] J. Brezovsky, E. Chovancova, A. Gora, A. Pavelka, L. Biedermannova, and J. Damborsky, "Software tools for identification, visualization and analysis of protein tunnels and channels," *Biotechnol*ogy Advances, vol. 31, no. 1, pp. 38 – 49, 2013.
- [12] M. Krone, B. Kozlíková, N. Lindow, M. Baaden, D. Baum, J. Parulek, H. C. Hege, and I. Viola, "Visual analysis of biomolecular cavities: State of the art," *Computer Graphics Forum*, vol. 35, no. 3, 2016.
- [13] D. Devaurs, L. Bouard, M. Vaisset, C. Zanon, I. Al-Bluwi, R. Iehl, T. Simon, and J. Corts, "MoMA-LigPath: a web server to simulate proteinligand unbinding," *Nucleic Acids Research*, vol. 41, no. W1, 2013.
- [14] R. Salomon-Ferrer, D. A. Case, and R. C. Walker, "An overview of the Amber biomolecular simulation package," Wiley Interdisciplinary Reviews: Computational Molecular Science, vol. 3, no. 2, pp. 198–210, 2013.
- [15] M. J. Abraham, T. Murtola, R. Schulz, S. Páll, J. C. Smith, B. Hess, and E. Lindahl, "GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers," *SoftwareX*, vol. 1-2, pp. 19 – 25, 2015.
- [16] D. Branduardi, F. L. Gervasio, and M. Parrinello, "From A to B in free energy space," *The Journal of Chemical Physics*, vol. 126, no. 5, p. 054103, 2007.
- [17] P. Tiwary, V. Limongelli, M. Salvalaglio, and M. Parrinello, "Kinetics of protein–ligand unbinding: Predicting pathways, rates, and rate-limiting steps," *Proceedings of the National Academy of Sciences*, vol. 112, no. 5, pp. E386–E391, 2015.
 [18] M. Suan Li and B. Khanh Mai, "Steered molecular dynamics-
- [18] M. Suan Li and B. Khanh Mai, "Steered molecular dynamicsa promising tool for drug design," *Current Bioinformatics*, vol. 7, no. 4, pp. 342–351, 2012.
- [19] P. H. Lee, K. L. Kuo, P. Chu, E. M. Liu, and J. H. Lin, "SLITHER: a web server for generating contiguous conformations of substrate molecules entering into deep active sites of proteins or migrating through channels in membrane transporters," *Nucleic Acid Research*, vol. 37, no. W559-64, 2009.
- [20] K. Fischer, "Smallest enclosing balls of balls," Ph.D. dissertation, Swiss Federal Institute of Technology, ETH Zurich, 2005.
- [21] S. J. Wright and J. Nocedal, Numerical Optimization. Springer, 1999.
- [22] S. M. Marques, D. Bednar, and J. Damborsky, "Computational study of protein-ligand unbinding for enzyme engineering," *Frontiers in Chemistry*, vol. 6, 2019.
- [23] D. B. Kokh, M. Amaral, J. Bomke, U. Griädler, D. Musil, H.-P. Buchstaller, M. K. Dreyer, M. Frech, M. Lowinski, F. Vallee, M. Bianciotto, A. Rak, and R. C. Wade, "Estimation of drug-target residence times by π-random acceleration molecular dynamics simulations," *Journal of Chemical Theory and Computation*, vol. 14, no. 7, pp. 3859–3869, 2018.



Jiří Filipovič holds B.Sc. and M.Sc. in Applied Informatics (Masaryk University) with specialization on numerical computing and Ph.D. in Informatics (Masaryk University). In 2012, he received the 1st prize in Joseph Fourier Award in Computer Science. After defending Ph.D., he worked as a postdoc at Masaryk University and the University of Vienna. Since 2017, he is head of High Performance Computing research group in CERIT-SC Centre at Institute of Computer Science, Masaryk University. His research inter-

ests include scientific and high performance computing, in particular methods for autotuning, source-to-source code transformation, heterogeneous computing, simulation, and computational biology.



Ondřej Vávra holds B.Sc. in Molecular biology and Genetics and M.Sc. in Genomics and Proteomics (Masaryk University). Currently he is working on Ph.D. in Ecotoxicology (Masaryk University). In 2017, he received the Brno Ph.D. Talent award. He is a member of Loschmidt Laboratories since 2013. His research interests include bioinformatics, structural and computational biology.



Jan Plhák holds B.Sc. in Mathematics (Masaryk University) and M.Sc. in Mathematics with Informatics (Masaryk University). His professional interests include GPU acceleration, geometrical algorithms, and routing in large graphs. Nowadays, he is head of research and development team developing a search engine in Kiwi.com.



David Bednář holds B.Sc. in Biochemistry, B.Sc and M.Sc. in Molecular biology and Genetics (Masaryk University), and Ph.D. in Molecular and Cellular Biology (Masaryk University) with specialization in molecular modelling in enzymology. In 2011, he received the Ph.D. talent award. After defending PhD, he worked as a postdoc at the University of North Carolina, Chapel Hill, USA. Since 2017, he is the assistant professor in Loschmidt Laboratories, Department of Experimental Biology, Masaryk Uni-

versity. His research interests include molecular modelling, protein engineering, and the development of computational tools for enzyme analysis and protein engineering.

FILIPOVIČ, VÁVRA et al.: CAVERDOCK: A NOVEL METHOD FOR THE FAST ANALYSIS OF LIGAND TRANSPORT



Sérgio M. Marques obtained his Ph.D. in Chemistry in 2007, from the Instituto Superior Tcnico (University of Lisbon, Portugal). During his first postdoc, conducted at the Instituto Superior Tcnico and the University of Pisa (Italy), he worked on the development of multifunctional agents against Alzheimer's disease and cancer. Between 2014 and 2016 he held a Marie Curie fellowship at the Masaryk University, in Czech Republic, where he focused his research on the in silico study of molecular tunnels and gates

in proteins and their engineering. He currently works as a research fellow at Masaryk University, where his main research interests are the computational modelling of enzymatic systems with the aim of improving their activity and stability towards biotechnological applications. He has co-authored 35 articles, 1 book chapter and 2 patents.



Jan Brezovský received his M.Sc. in Biophysics at Masaryk University, Brno, Czech Republic in 2006. In 2011 he received Ph.D. in Environmental Chemistry with Thesis "Molecular Modeling of Enzymes" at the same university. Since then, he has been working as junior researcher at Loschmidt Laboratories of the Masaryk University focusing on computational protein engineering and molecular modeling. His main research topics are ration redesign of dynamics and selectivity of enzymes. Jan Brezovsky is currently

leading the Laboratory of Biomolecular Interactions and Transport, Department of Gene Expression, Institute of Molecular Biology and Biotechnology Faculty of Biology, Adam Mickiewicz University in Poland.



Luděk Matyska is a full professor of Informatics and since 2013 a director of Institute of Computer Science at Masaryk University; he also holds a senior researcher position at CES-NET. He leads national research e-infrastructure CERIT-SC and is active internationally in grid, cloud and recently EOSC-related activities; he holds a position of EGI Executive Board member and is also EOSC-Hub Project Management Board member. His research interest lies in security in large-scale distributed systems, e-

infrastructure (network, computing and data) architecture and policy and generally in cloud and grid systems. He authored or co-authored more than 100 papers and conference contributions.



Jiří Damborský is the Josef Loschmidt Chair Professor of Chemistry at the Faculty of Science at Masaryk University and a group leader at the International Centre for Clinical Research in Brno, Czech Republic. Research of his group focuses on protein engineering. His group develops new concepts and software tools for protein engineering (CAVER, CAVERDOCK, HOTSPOT WIZARD, PREDICTSNP, FIREPROT, CALFIT-TER, SOLUPROT, ENZYMEMINER), and uses them for the rational design of enzymes. He has

published 200 original articles, 15 book chapters and filed 6 international patents. He is a co-founder of the first biotechnology spin-off from Masaryk University Enantis Ltd. Among the awards and distinctions he has received is the award EMBO/HHMI Scientist of the European Molecular Biology Organisation and the Howard Hughes Medical Institute.