

HIGH-PERFORMANCE COMPUTING ON BIOINFORMATICS

Anton Stoilov, Borislav Yurukov

1. INTRODUCTION

High-performance computing (HPC) is an important domain of the computer science field. For more than 30 years, it has allowed finding solutions to problems and enhanced progress in many scientific areas such as bioinformatics and drug design. The binding of small molecule ligands to large protein targets is central to numerous biological processes. The accurate prediction of the binding modes between the ligand and protein (the docking problem) is of fundamental importance in modern structure-based drug design. The interactions between the receptor and ligand are quantum mechanical in nature, but due to the complexity of biological systems, quantum theory cannot be applied directly. Consequently, most methods used in docking and computational drug discovery are more empirical in nature and usually lack generality. Quantum mechanical phenomena, such as the formation of a covalent bond between the protein and the ligand upon binding during the transition state of the reaction, cannot be predicted and/or evaluated using these empirical methods. In the field of molecular modeling, docking is a method which predicts the preferred orientation of one molecule to a second when bound to each other to form a stable complex. Knowledge of the preferred orientation in turn may be used to predict the strength of association or binding affinity between two molecules using, for example, scoring functions. Docking is frequently used to predict the binding orientation of small molecule drug candidates to their protein targets in order to in turn predict the affinity and activity of the small molecule. Hence docking plays an important role in the rational design of drugs. Given the biological and pharmaceutical significance of molecular docking, considerable efforts have been directed

towards improving the methods used to predict docking. Each docking program makes use of one or more specific search algorithms, which are the methods used to predict the possible conformations of a binary complex. An overview of current docking techniques is presented with a description of applications including a benchmark for docking on IBM HPC platform, also mathematical algorithm will be presented. The present benchmark is made from an existing test set (CCDC/Astex Validation Set) on typical HPC system. Selected examples were docked with GOLD software.

2. MOLECULAR DOCKING

Molecular docking is a computer simulation procedure to predict the conformation of a receptor-ligand complex, where the receptor is usually a protein or a nucleic acid molecule (DNA or RNA) and the ligand is either a small molecule or another protein. It can also be defined as a simulation process where a ligand position is estimated in a predicted or pre-defined binding site. Molecular docking research focusses on computationally simulating the molecular recognition process. It aims to achieve an optimized conformation for both the protein and ligand and relative orientation between protein and ligand such that the free energy of the overall system is minimized.

Computational docking of a small molecule to a biological target involves efficient sampling of possible poses of the former in the specified binding pocket of the latter in order to identify the optimal binding geometry, as measured by a user-defined fitness or score function. X-ray crystallography and NMR spectroscopy continue to be the primary source of 3-dimensional structural data for protein and nucleic acid targets. In favorable cases where proteins of unknown structure have high sequence homology to known structures, homology modeling can provide a viable alternative by generating a suitable starting point for “in silico” discovery of high affinity ligands. Potential energy of molecular field model is a function of a atomic position (x,y,z) normally in Cartesian space. The equation of the potential energy of the system of atoms in the molecular force field, commonly used in molecular modeling is presented below:

$$\begin{aligned}
 (1) \quad E = & \frac{1}{2} \sum_{i=1}^m k_{\theta_i} (\theta_i - \theta_{0,i})^2 + \frac{1}{2} \sum_{i=1}^n k_{b_i} (b_i - b_{0,i})^2 + \\
 & \frac{1}{2} \sum_{i=1}^k v_i [1 + \cos(n_i \omega_i - \gamma_i)] + \\
 & \sum_{i,j} 4\epsilon_{i,j} \left[\left(\frac{\sigma_i}{r_{ij}} \right)^{12} - \left(\frac{\sigma_j}{r_{ij}} \right)^6 \right] + \sum_{i,j} \frac{1}{4\pi\epsilon_0} \frac{q_i q_j}{r_{ij}^2}
 \end{aligned}$$

The complexity of computational docking increases in the following order:

(a) rigid body docking, where both the receptor and small molecule are treated as rigid. (b) flexible ligand docking, where the receptor is held rigid, but the ligand is treated as flexible; (c) flexible docking, where both receptor and ligand flexibility is considered.

Docking applications usually make use of one or more of the following methods: fast shape matching (SM), incremental construction (IC), Monte Carlo simulations (MC), distance geometry (DG), evolutionary programming (EP), genetic algorithms (GA), tabu search (TS) and simulated annealing (SA). The GA adopted by GOLD algorithm requires as input the approximate size and location of the receptor active site and also the coordinates of protein and a ligand conformation. The active site may be defined by several techniques. GA is also implemented in the program DOCK, which is able to dock either whole ligand inside active site or a rigid fragment of the ligand. "Lamarckian" GA (LGA) is also implemented in docking algorithms. The LGA switches between "genotypic space" and "phenotypic space." Mutation and crossover occur in genotypic space, while phenotypic space is determined by the energy function to be optimized. Energy minimization (local sampling) is performed after genotypic changes have been made to the population (global sampling) in phenotypic space, which is conceptually similar to MC minimization.

3. METHODS

GOLD software uses a Genetic Algorithm (GA) for protein ligand docking which works as follows:

- A population of potential solutions (i.e. possible docked orientations of the ligand) is set up at random.
- Each member of the population is encoded as a chromosome, which contains information about the mapping of protein ligand interactions.
- Each chromosome is assigned a fitness score based on its predicted binding affinity and the chromosomes within the population are ranked according to fitness.

The population of chromosomes is iteratively optimised. At each step, a point mutation may occur in a chromosome, or two chromosomes may mate to give a child. The selection of parent chromosomes is biased towards fitter members of the population, i.e. chromosomes corresponding to ligand dockings with good fitness scores. The GOLD validation test set is one of the most comprehensive of all of the docking methods reviewed, and achieved a 71% success rate based primarily on a visual inspection of the docked structures. 66 of the complexes had an RMSD of 2.0 Å or less, while 71 had an RMSD of 3.0 Å or less. The omission of hydrophobic interactions and a solvent model may explain some of the docking failures which included highly flexible, hydrophobic ligands, and those complexes containing poorly resolved active sites. However, recent extensions to GOLD include the addition of hydrophobic fitting points that are used in the least squares fitting algorithm to generate the ligand orientation.

In this paper the benchmark test set is based on CCDC/Astex Validation Set developed by Cambridge Crystallographic Data Centre (CCDC) for docking software GOLD. There are 305 entries and protonation states have been set in all cases. The CCDC/Astex set consists of the following files for each entire: protein.mol2 file

ligand_reference.mol2 file - this contains the ligand pose as found in the PDB entry. Entries with multiple binding modes, such as 1abe, are stored as follows: ligand_reference1.mol2, ligand_reference2.mol2, with the accompanying protein files protein1.mol2 and protein2.mol2.

gold.conf file - the GOLD configuration file can be used with the GOLD docking program. It also contains the centre and radius of the binding site. For covalently-bound ligands, a flag is set in this file and atom numbers of the link are stored.

water.mol2 file - this file is available for those PDB entries that include a water set; it is currently only available for entries that were not included in the previous GOLD validation set.

The equipment for experiments was provided by CENTER FOR ADVANCED BIOINFORMATICS RESEARCH, South-West University “Neofit Rilski”, Blagoevgrad, BULGARIA. This equipment include two different computational power with this characteristics:

IBM x3650 M2 – Processor: 2 x Xeon Quad-Core Intel Xeon 4C Model E5520 4C 2.26GHz with EM64T/1066MHz, 8MB L3 Cache, RAM: 2 x2GB 4096MB ECC PC3-10600 DDR3, HDD: 2x146GB 10K SAS Hot Swap RAID controller - Integrated RAID MR10i

IBM BladeCenter HS22 – Processor: 2 x Xeon Quad-Core Intel Xeon 4C Model E5504 80W 2.00GHz/800MHz/4MB L2, 2x2GB, O/Bay 2.5in SAS 1, RAM: 4 x 2GB Single Rank PC3-10600 CL9 ECC DDR3, HDD: 2 x IBM 146 GB 2.5in SFF Slim-HS 10K 6Gbps SAS HDD

4. RESULTS

After successful docking procedure for each test-case is observed by Fitness (scoring function), Best ranking time and Total run time . The results from experiments was present on Table 1. In the fields of molecular modelling, scoring functions are fast approximate mathematical methods used to predict the strength of the non-covalent interaction (also referred to as binding affinity) between two molecules after they have been docked. Most commonly one of the molecules is a small organic compound such as a drug and the second is the drug's biological target such as a protein receptor.

Table 1. 22 specific cases from CCDC/ASTEX validation test

Test ID	ligand byte	Fitness	IBM x3650 M2		IBM BladeCenter HS22	
			Bestranking time,s	Total run time,s	Bestranking time,s	Total run time,s
1a0q	982	88.82	34.18	37.30	32.14	36.23
1a1b	1430	100.51	196.37	203.20	194.80	198.12
1a1e	1430	99.40	184.43	191.06	182.20	190.47
1a4g	1073	91.52	36.02	40.46	32.10	35.38
1a4k	985	65.61	116.88	122.45	113.55	117.47
1a4q	1074	102.89	70.83	77.14	65.32	67.17
1a6w	807	60.22	21.53	22.55	20.46	21.57
1a07	1432	53.44	173.49	176.88	170.21	174.22
1a9u	1163	64.29	30.89	35.67	28.23	33.76
1a28	988	68.03	34.89	36.05	32.90	35.54
1a42	893	77.93	67.90	70.02	62.56	66.89
1aaq	545	94.53	259.27	263.46	249.18	253.38

Test ID	ligand byte	Fitness	IBM x3650 M2		IBM BladeCenter HS22	
			Bestranking time,s	Total run time,s	Bestranking time,s	Total run time,s
1abe	542	54.22	20.24	41.16	18.31	40.03
1abe	542	57.69	20.19	41.16	18.15	40.03
1abf	542	58.47	23.36	47.62	21.35	45.62
1abf	542	55.19	23.45	47.62	21.59	45.62
1acj	537	52.84	20.50	21.45	17.85	20.58
1acl	634	75.33	99.15	100.15	96.23	99.37
1acm	901	96.64	22.05	23.03	20.37	23.01
1aco	634	86.49	17.51	17.97	13.05	15.85
1aec	812	63.00	49.23	50.01	45.28	48.68
1aha	545	47.15	13.47	14.08	11.42	13.27
1ai5	988	47.59	15.93	16.83	13.25	15.49
1aj7	988	82.48	37.91	40.13	33.06	38.86

5. CONCLUSIONS

An extensive summary of currently available docking methods has been presented. Comparisons suggest that the best algorithm for docking is probably a hybrid of various types of algorithm encompassing novel search and scoring strategies. The most useful docking method will not only perform well, but will be easy to use and parametrise, and sufficiently adaptable such that different functionality may be selected, depending on the number of structures to be docked, the available computational resources, and the complexity of the problem. If the parameters cannot be generated quickly then although the algorithm may be computationally efficient, from a practical point of view it is limited. Conversely, a rapid scoring function may not necessarily be able to model some specific interactions. Moreover, although current docking methods show great promise, fast and accurate discrimination between different ligands based on binding affinity, once the binding mode is generated, is still a significant problem.

ACKNOWLEDGMENT

This research was supported by projects:

- Bioinformatics Research: Protein Folding, Docking and Prediction of Biological Activity NSF - 02/16 funded by National Research Fund of Bulgaria and South-West University "Neofit Rilski".
- Bioinformatical studies on the structure and activity of proteins and drug – receptor interactions (BIOINFO)
- Interregional Cooperation at Scientific Computing in Interdisciplinary Science (ICOSCIS)

REFERENCES

1. Halperin I, Ma B, Wolfson H, Nussinov R (June 2002). "Principles of docking: An overview of search algorithms and a guide to scoring functions". *Proteins* 47 (4): 409–443.
doi:10.1002/prot.10115. PMID 12001221.
2. Mustard D, Ritchie DW (August 2005). "Docking essential dynamics eigenstructures". *Proteins* 60 (2): 269–274.
doi:10.1002/prot.20569. PMID 15981272.
3. Shoichet BK, Stroud RM, Santi DV, Kuntz ID, Perry KM (March 1993). "Structure-based discovery of inhibitors of thymidylate synthase". *Science* 259 (5100): 1445–50.
doi:10.1126/science.8451640. PMID 8451640.
4. McGann MR, Almond HR, Nicholls A, Grant JA, Brown FK (January 2003). "Gaussian docking functions". *Biopolymers* 68 (1): 76–90.
doi:10.1002/bip.10207. PMID 12579581.
5. Friesner RA, Banks JL, Murphy RB, Halgren TA, Klicic JJ, Mainz DT, Repasky MP, Knoll EH, Shelley M, Perry JK, Shaw DE, Francis P, Shenkin PS (March 2004). "Glide: a new approach for rapid, accurate docking and scoring. 1. Method and assessment of docking accuracy". *J. Med. Chem.* 47 (7): 1739–1749.
doi:10.1021/jm0306430. PMID 15027865.
6. Jain AN (February 2003). "Surflex: fully automatic flexible molecular docking using a molecular similarity-based search engine". *J. Med. Chem.* 46 (4): 499–511. doi:10.1021/jm020406h. PMID 12570372.
7. Zsoldos Z, Reid D, Simon A, Sadjad SB, Johnson AP (July 2007). "eHiTS: a new fast, exhaustive flexible ligand docking system". *J. Mol. Graph. Model.* 26 (1): 198–212.
doi:10.1016/j.jmglm.2006.06.002. PMID 16860582.
8. Jones G, Willett P, Glen RC, Leach AR, Taylor R (April 1997). "Development and validation of a genetic algorithm for flexible docking". *J. Mol. Biol.* 267 (3): 727–748.
doi:10.1006/jmbi.1996.0897. PMID 9126849.

◇ **Anton Stoilov**

South-West University "Neofit Rilski",
Faculty of Engineering, Department of Electrical Engineering,
Electronics and Automatics

◇ **Borislav Yurukov**

South-West University "Neofit Rilski",
Faculty of Mathematics and Natural Sciences,
Department of Informatics

