

DOCKING: **METHODS OVERVIEW**

VLAD KHOLODOVYCH

**UNIVERSITY OF MEDICINE & DENTISTRY
OF NEW JERSEY (UMDNJ)**

*KHOLODVL@UMDNJ.EDU
HTTP://WWW2.UMDNJ.EDU/~KHOLODVL*

DOCKING

WHAT:

Prediction of the optimal configuration and energy between two molecules

WHY:

Assumption: components that dock well may bind to each other

HOW:

Change orientations of molecules to maximize their interaction

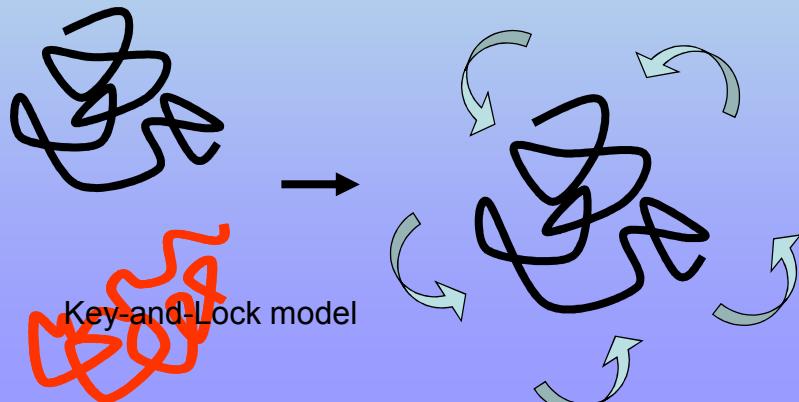
Minimize the total energy of interaction
(*The best binding pose has the minimal binding energy*)

RIGID DOCKING

Both molecules are rigid

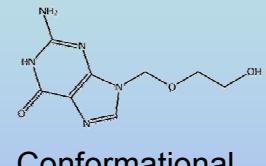
Orientation change: translation, rotation

No change in conformation

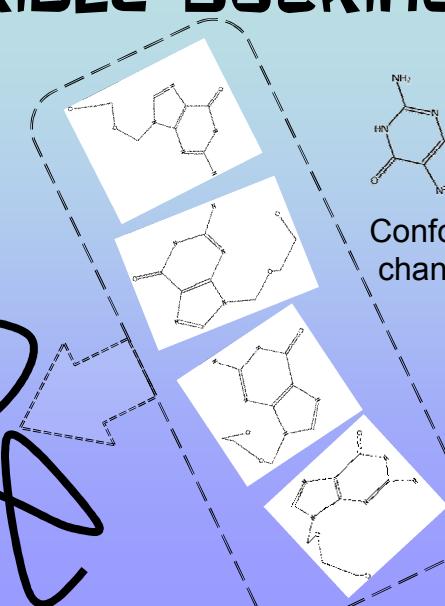
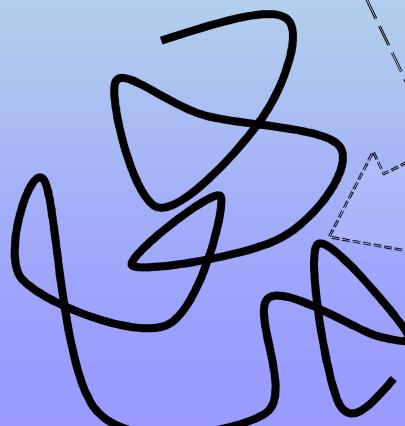


FLEXIBLE DOCKING

Rigid receptor



Conformational changes in ligand



TYPE OF DOCKING SIZE WISE

PROTEIN-PROTEIN DOCKING

- Bound docking (“rigid” redocking):
6 degrees of freedom: 3 for rotation, 3 for translation
- Unbound docking : side chain flexibility

PROTEIN-SMALL MOLECULE DOCKING

- Rigid receptor, rigid ligand
- **Rigid receptor, flexible ligand**
- Flexible receptor, flexible ligand

RECEPTOR - LIGAND DOCKING: SPECIFICS

- Many possible sites for potential interactions
- Extremely large search space
 - Both molecules are flexible – hundreds to thousands of degrees of freedom (DOF)
 - Enormous number of possible docking poses

Receptor - Ligand Docking

Known protein.
No ligand, no information about site interactions

Receptor - Ligand Docking

Known protein.
No ligand, no information about site interactions

CASTp Computed Atlas of Surface Topography of proteins

<http://sts-fw.bioengr.uic.edu/castp/calculation.php>

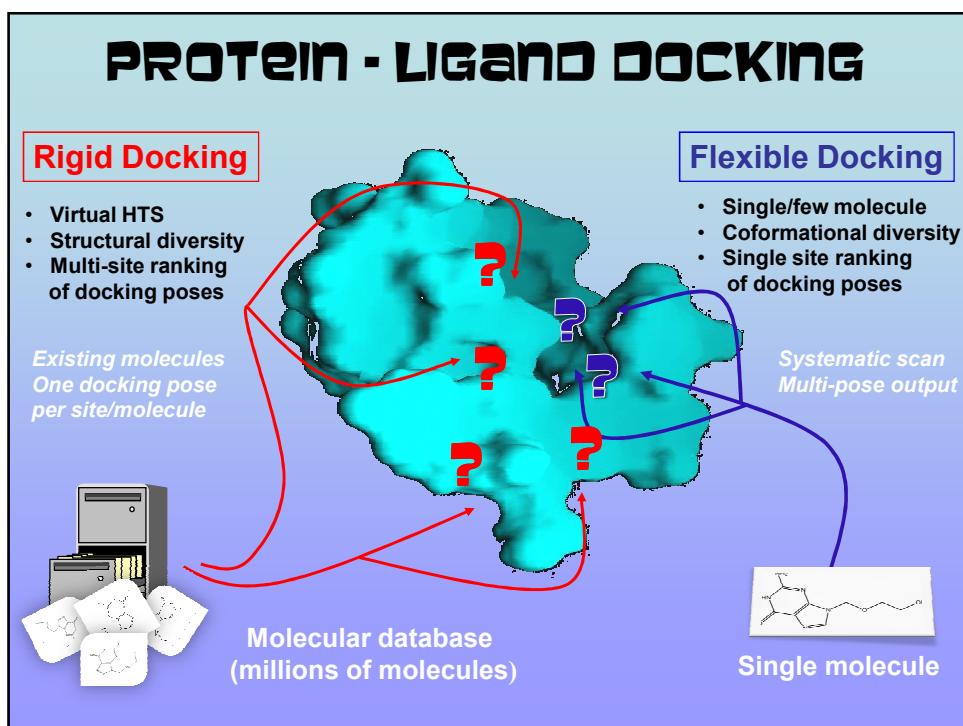
ID	Area	Vol
32	964.4	1984.8
21	178.0	177.4
20	118.7	174.6
28	135.7	149.8
29	88	117.0
27	81.5	102.0
26	41.7	61
25	30.2	82.2
24	34	33.8
23	06	66.7
22	26.2	26.4
21	98.8	47.9
20	16.1	29.1

Please cite CASTp as Liang J. et al. Protein Sciene, 1998, 7:1884¹

Pocket color: pocket02 [green] Display: Wireframe [x] Enter RasMol commands below: Quick Reference Run

Open in: [180](#) [SEEEAKAWKAQVWVHIIKPLIIVKIGIDRERRGAVWWVVKPAPAKREKEFELIIP](#)
[242](#) [HHAADMSTYLMVLSIPLAKVSYKHLHIEEYSLVAAAKKLGCKLHIV](#)
[250](#) [TACTCTTWERCGDLYKQKQKACGCGTGTGAGCAGGTCAGKPKHUYMLKKGLOLCEYY](#)

Show annotated sites in Red

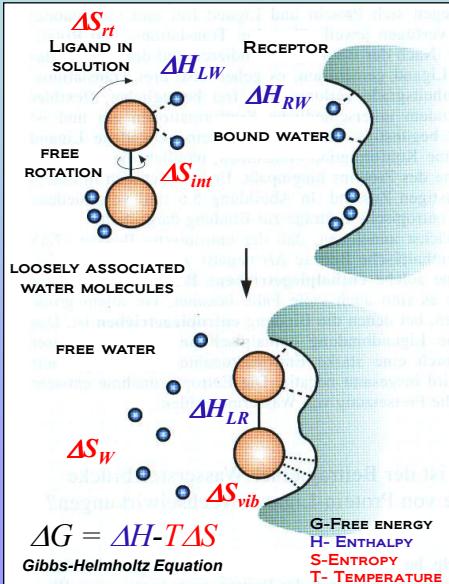


KNOWN TARGET PROTEIN: SCREENING

	Rigid	Flexible
Database search	YES	NO
Multi-conformational	NO	YES
Basis	<i>Volume/Surface complementarity</i>	<i>Interaction Energy</i>
Algorithm	<i>Clique matching (shape)</i>	<i>Simulated Annealing</i> <i>Genetic Algorithm</i> <i>Fragment Building</i>
Programs	<i>Dock</i> <i>GOLD</i> <i>GRAMM</i> <i>MOE</i>	<i>AutoDock, Dock</i> <i>FlexX</i> <i>GOLD</i> <i>MOE</i>

DOCKING: DIFFICULTIES

RECEPTOR-LIGAND COMPLEX



Multiple steps in the receptor – ligand interaction:

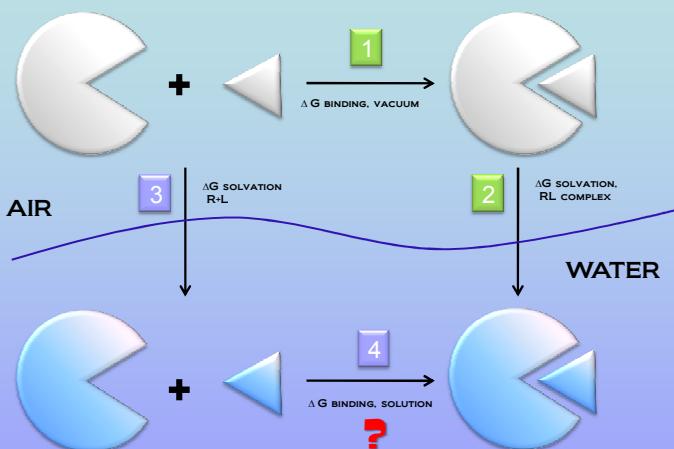
- Approach
- Desolvation of the ligand and the binding site of a protein
- Penetration into the protein cavity
- Change of the ligand orientation
- Adoption of the correct “active” conformation
- Establishing of new H-bonds, electrostatic and hydrophobic contacts

In general free energy function is formed from several terms

$$\Delta G_{bind} = \Delta G_{vdw} + \Delta G_{hbond} + \Delta G_{elect} + \Delta G_{conform} + \Delta G_{tor} + \Delta G_{sol}$$



SOLVATION BINDING



Hess's law of heat summation:

The change in free energy between two states will be the same, no matter what the path.

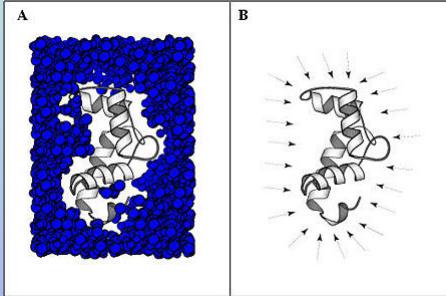
Tricky math
 $1+2 = 3+4$
 $4 = 1+2-3$

The energy of binding in solvent can be calculated as followed:

$$\Delta G_{binding,solution} = \Delta G_{binding,vacuum} + \Delta G_{solvation}(RL) - \Delta G_{solvation}(R+L)$$

4 1 2 3

SOLVATION: EXPLICIT VS IMPLICIT



the N terminal domain of Troponin

A A view of the inside of the initial box for a simulation which includes water molecules explicitly; the whole system consists of **1229** atoms in the protein and **6009** atoms for water.

B In a implicit solvent simulation, the effect of water is included in a potential of mean force, $W(X)$, visualized as arrows. In this case, the simulation only includes the 1229 protein atoms.

Generalized Born / Surface Area - GB/SA model

$$G_{solv} = G_{cav} + G_{vdW} + G_{pol}$$

the total solvation free energy G_{solv} is given as the sum of a solvent-solvent term (G_{cav}), a solute-solvent van der Waals term (G_{vdW}), and a solute-solvent electrostatics polarization (G_{pol})

$$G_{cav} + G_{vdW} = \sum_{k=1}^N ASP(k) ASA(k)$$

$G_{cav} + G_{vdW}$ are computed together as a function of the solvent-accessible surface areas. where ASA(k) is the total solvent accessible surface area of atom k and ASP(k) is an empirical atomic solvation parameter.

$$G_{pol} = -166.0 \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \sum_{j=1}^N \frac{q_i q_j}{\sqrt{r_{ij}^2 + b_i b_j \exp\left(-\frac{r_{ij}^2}{4b_i b_j}\right)}}$$

$$G_{Born} = -\frac{q^2}{2b} \left(1 - \frac{1}{\epsilon}\right)$$

Born Equation

where the double sum runs over all pairs of atoms (i and j). q_i and q_j are the partial charges of i and j, respectively, and r_{ij} is the i,jth atom pair separation. ϵ is the dielectric constant of the medium. b_i and b_j are the so-called Born radius of atom i and j, respectively .

Born radius b is the effective radius of the atom defined by Born equation

TAKE HOME MESSAGE

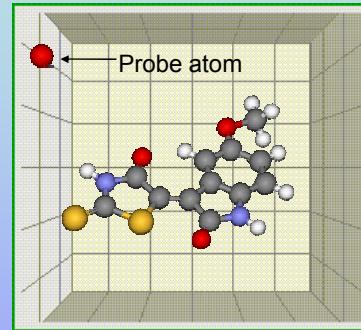
- **Binding energy** is a multicomponent function that is not easy to calculate
- Solvation component of binding energy is computationally demanding to be solved **explicitly**
- Instead, the **implicit** solvation term is included in a potential of mean force field, while only essential solvent molecules mediated protein-ligand interactions in the binding pocket are included explicitly
- GB/SA method is the most popular way to calculate solvation
- Docking programs use simple and fast methods to estimate binding interactions. Solvation factor is not included in the scoring function.

TWO IMPORTANT QUESTIONS IN DOCKING

1. How to align (place) the ligand molecule into the binding pocket?
2. How to alter a ligand structure?

DOCKING ALIGNMENT: GRID

- Grid-based docking (AutoDock)
- Ligand-protein interaction energies are pre-calculated and then used as a look-up table during simulation

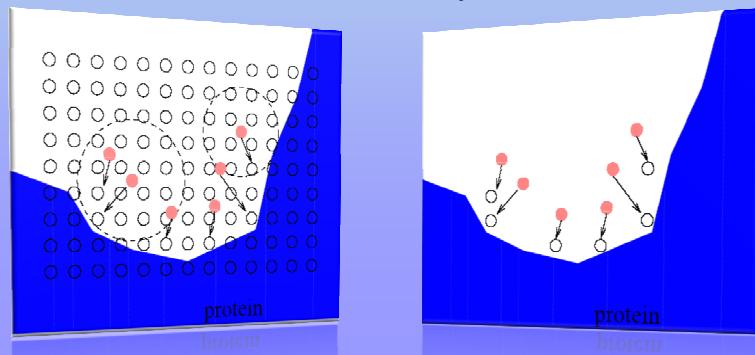


- Grid maps are constructed for each atom type presented in the ligand

DOCKING ALIGNMENT: GRID

Each ligand atom is matched with the grid point with the lowest energy within its neighborhood

Total binding energy is estimated by summation from different “element” maps



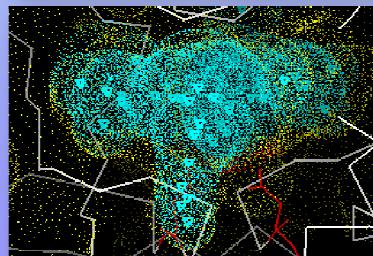
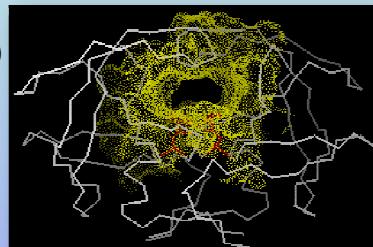
DOCKING ALIGNMENT : SPHERES

Spheres-based docking (UCSF Dock)

The receptor cavities are filled up with spheres.

Sphere centers are then matched to the ligand atoms to determine possible orientations for the ligand. Typically 10^4 - 10^5 orientations are generated for each ligand molecule.

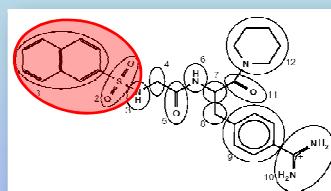
The sphere centers are identified by cyan triangles, and the sphere surfaces are shown in dots



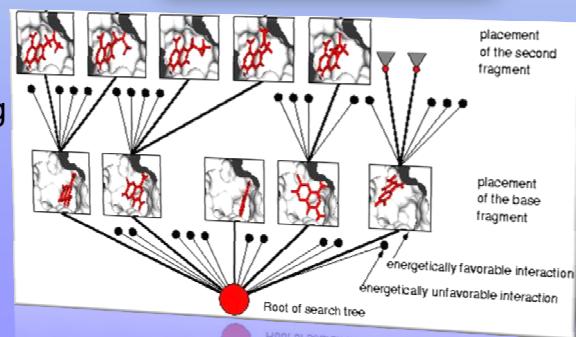
DOCKING ALIGNMENT: FRAGMENT ADDITION

FlexX and DOCK

- Split the ligand into fragments
- Determine the rigid core/root



- Incrementally add the fragments to the core/root frame during docking optimization



LIGAND MODIFICATION ALGORITHMS

Simulated Annealing

- Based on temperature effects
- Start with high temperature —→ Global search
- Lower temperature —→ Local search

Genetic Algorithm

- Charles Darwin's Theory of Evolution
Genotype —→ Phenotype
- Lamarckian Algorithm (Jean –Baptiste de Lamarck)
Phenotype —→ Genotype

SIMULATED anneALING

Simulated annealing is a global optimization technique based on the Monte Carlo method.

Generate small random changes in the current state (ligand coordinates)

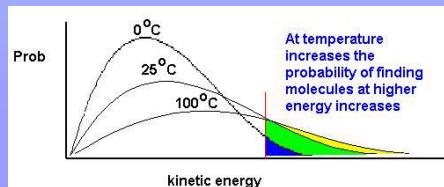
Metropolis criterion :

moves that decrease the energy of the system are always accepted
moves that increase the energy of the system are accepted only
if probability p is larger than random chosen number between 0 and 1

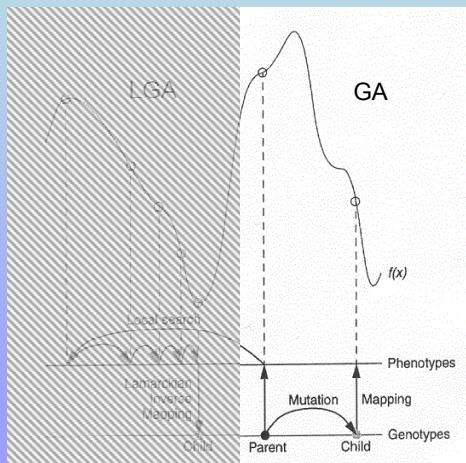
$$p = e^{\left(\frac{-\Delta E}{kT}\right)} \quad \text{where } \Delta E = E_1 - E_0$$

E_1 is an energy of the new state, E_0 is an energy of the current state,
 T is the "temperature" of the simulation, and k is Boltzmann's constant

At high temperatures, many states are accepted, while at low temperatures, the majority of these probabilistic moves are rejected



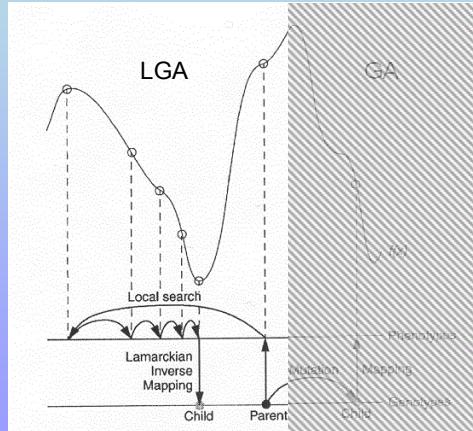
GENETIC ALGORITHM



- In standard GA, the genotype (x, y, z coordinates plus rotation and any torsion angles) are mapped to the fitness function $f(x)$
- All new generation based only on parents genes
- Genotypes of parents with high $f(x)$ values are mutated forming genotypes of children with lower $f(x)$ values

LAMARCKIAN GENETIC ALGORITHM

- LGA finds lowest fitness function (energy) values first, then maps these values to their respective genotypes
- Each new child is allowed to create a new generation
- Genetic algorithm plus Solis and Wets local search
- Better performance than either simulated annealing or genetic algorithm alone



SCORING FUNCTIONS

Purpose: Evaluate the placement (and possibly a conformation) of the ligand-protein pair.

FORCE FIELD FUNCTIONS: use non-bonded energies of force fields (FlexiDock, Tripos, Amber, CHARMM)

EMPIRICAL SCORING FUNCTIONS: derived from a set of protein-ligand complexes with measured binding affinity (more chemically appealing, however require significantly more statistical fitting than force fields functions (ChemScore, AutoDock))

KNOWLEDGE-BASED SCORING FUNCTIONS: a statistical analysis of structures of protein-ligand complexes derived from databases, e.g. PDB. ("PMF", DrugScore).

Each scoring function has been derived from a different set of crystal structures and experimental data. Thus, it is reasonable to use multiple functions when evaluating docking results.

Consensus: structures which are considered as a good fit by multiple scoring functions should be tested first

C-SCORE

- **F-SCORE** FlexX scoring function (ref: Rarey, M.; Kramer, B.; Lengauer, T., Klebe, G. J. *J. Mol. Biol.* (1996) **261**, 470-489)
- **G-SCORE** Includes terms for hydrogen bonding, complex (ligand-protein) and internal (ligand-ligand) energies (ref : Jones, G. Willett, P., Glen, R. Leach, A.R., and Taylor, R. *J. Mol. Biol.* (1997) **267**, 727)
- **PMF-SCORE** Potential Mean Force Developed Helmholtz free energies of interactions for protein-ligand atom pairs (ref: Muegge, I., Martin, Y.C., *J. Med. Chem.* (1999) **42**, 791)
- **D-SCORE** Uses only charge and van der Waals interactions between the protein and the ligand (ref: Kuntz, I.D., Blaney, J.M., Oatley S.J., Langridge, R., Ferrin, T.E., *J. Mol. Biol.* (1982) **161**, 269)
- **CHEMSCORE** Includes terms for H-Bonding, metal-ligand interaction, lipophilic contact, and intercept term and rotational entropy (ref: Eldredge, M.D., Murray, C.W., Auton, T.R., Paolini, G.V., and Mee, R.P. *J. Comp.-Aided Molec Des.*,(1997) **11**, 425)
- **C-SCORE** is the sum of "good" score marks, min 0 to max 5.

PROGRams

PROTEIN-PROTEIN DOCKING

- **FTDock** algorithm maps both molecules onto orthogonal grids and performs a global scan of translational and rotational space. The scoring method is primarily a surface complementarity score between two grids.

<http://www.bmm.icnet.uk/docking/ftdock.html>

- **MULTIDock** is a program for refinement of FTDock results. Allows generation of side chain conformations based on a rotamer libraries.

www.bmm.icnet.uk/docking/multidock.html

PROTEIN-PROTEIN DOCKING

GRAMM:

GLOBAL RANGE MOLECULAR MATCHING

Main features of GRAMM:

- requires only the atomic coordinates of the molecules (no information about the binding sites is needed)
- performs an exhaustive 6-dimensional search through the relative translations and rotations of the molecules
- molecular pairs: two proteins, a protein and a smaller compound, two transmembrane helices, etc.

http://vakser.bioinformatics.ku.edu/main/resources_gramm1.03.php

LIGAND-PROTEIN DOCKING

NAME: AutoDock

AUTHORS: Olson et al 1990, Morris et al 1998

PUBLISHER: Scripps

LIGAND ALIGNMENT:

Grid

SEARCH ALGORITHMS:

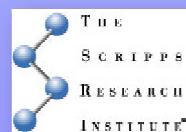
Simulated Annealing

Genetic Algorithm

Lamarckian GA (GA+LS hybrid)

UNIQUE FEATURES: Solvation term included in the scoring function, LGA

<http://autodock.scripps.edu>



LIGAND-PROTEIN DOCKING

NAME: UCSF DOCK

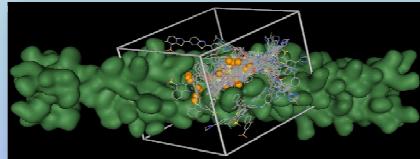
AUTHORS: Kuntz et al. 1982;
Ewing & Kuntz 2001

PUBLISHER: UCSF

LIGAND ALIGNMENT:

Spheres

SEARCH ALGORITHMS:
Fragment-Based



UNIQUE FEATURES: Ligand is matched to spheres, grid scoring

<http://dock.compbio.ucsf.edu/>

LIGAND-PROTEIN DOCKING

NAME: ICM Dock

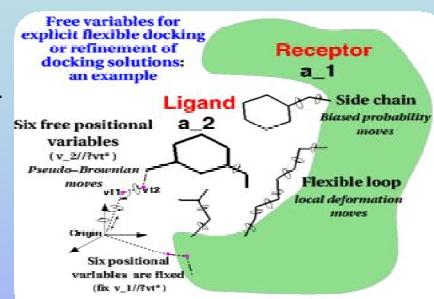
AUTHORS: Abagyan et al. 1994

PUBLISHER: MolSoft

LIGAND ALIGNMENT:

Grid

SEARCH ALGORITHMS:
Genetic Algorithm



UNIQUE FEATURES: Protein/Protein, Protein/Peptide,
Protein/Ligand docking; Partial flexibility of the protein

<http://www.molsoft.com/docking.html>



LIGAND-PROTEIN DOCKING



Cambridge Crystallographic Data Centre

NAME: GOLD

AUTHORS: Jones et al 1997

PUBLISHER: Cambridge Crystallographic Data Center

LIGAND ALIGNMENT:

Fitting Points

SEARCH ALGORITHMS:

Genetic Algorithm

UNIQUE FEATURES: Automatic consideration of cavity bound water molecules; Partial protein flexibility (side chains and backbone for up to 10 residues)

Free on-line test run at:

http://gold.ccdc.cam.ac.uk/setup_and_run_docking.php

<http://www.ccdc.cam.ac.uk/prods/gold/index.htm>

LIGAND-PROTEIN DOCKING

NAME: MOE Dock

AUTHORS:

PUBLISHER: Chemical Computing Group

LIGAND ALIGNMENT:

Spheres, pharmacophore triad

SEARCH ALGORITHMS:

Simulated Annealing

Genetic Algorithm

UNIQUE FEATURES: Pharmacophore mapping/filtering, solvation term is included in ligand conformers generation



<http://www.chemcomp.com/software-sbd.htm>

LIGAND-PROTEIN DOCKING

NAME: FlexX

AUTHORS: Rarey, Lengauer 1996

PUBLISHER: BioSolvEit

LIGAND ALIGNMENT:

Fragment Based

SEARCH ALGORITHMS:

Fragment-Based



UNIQUE FEATURES: Pharmacophore based docking; docking into ensemble of protein structures; automatically assigned metal coordination geometry; high speed vHTS

<http://www.biosolveit.de/FlexX/>

LIGAND-PROTEIN DOCKING

NAME: GLIDE



AUTHORS: Friesner et al 2004

PUBLISHER: Schrodinger

LIGAND ALIGNMENT:

Grid

SEARCH ALGORITHMS:

Simulated Annealing

UNIQUE FEATURES:

<http://www.schrodinger.com/ProductDescription.php>

EVALUATION OF DOCKING

ROOT MEAN SQUARE
DEVIATION (*RMSD*)

$$RMSD = \sqrt{\frac{\sum d_i^2}{n}}$$

n - number of atoms

d_i - distance between two corresponding atoms *i*
in two structures

UNIT OF RMSD: ÅNGSTROMS

IDENTICAL STRUCTURES: *RMSD* = "0"

SIMILAR STRUCTURES: *RMSD* IS SMALL (1 – 2 Å)

DISTANT STRUCTURES : *RMSD* > 3 Å

moe and GOLD DOCKING TUTORIAL

[HTTP://WWW2.UMDNJ.EDU/~KHOLODVL](http://WWW2.UMDNJ.EDU/~KHOLODVL)