

# A Virtual Screening Procedure Combining Pharmacophore Filtering and Molecular Docking with the LIE method

by

Guzin Tunca

A Thesis Submitted to  
Universitat Autònoma de Barcelona  
in Partial Fulfillment of the Requirements for  
the Degree of  
Doctorate of Philosophy

Tesi Doctoral UAB /ANY 2012

Director de la tesi: Xavier Daura

Institut de Biotecnologia i de Biomedicina



This is to certify that I have examined this copy of a doctoral dissertation by

Guzin Tunca

and have found that it is complete and satisfactory in all respects,

and that any and all revisions required by the final

examining committee have been made.

**Committee Members:**

Dr Leonardo Pardo

---

Dr. Amadeu Llebaria

---

Dr. Jean-Didier Maréchal

---

**Thesis Supervisor:**

Dr. Xavier Daura

---



*To my family*



## *Acknowledgements*

I would like to thank to my thesis Supervisor Dr. Xavier Daura for his advice and support during the five years. I am also grateful to Dr. Roman Affentranger, who was the best post-doc in the world for his great help at every step of the computational workflow created in this study, this study wouldn't be possible without him. Dr. Hugo Gutierrez de Teran gave invaluable information and help about the Linear Interaction Energy method and I am thankful to him for providing the necessary scripts and detailed explanation of the LIE technique. Dr. Xavier Barril from University of Barcelona kindly let us use the small molecules library created and curated by his group. My thanks also go to Dr. Amadeu Llebaria and his group from CSIC for carrying out the experimental testing of the ligand candidates for glucocerebrosidase. I also thank to Dr. Leonardo Pardo, Dr. Jean-Didier Marechal and Dr. Amadeu Llebaria for accepting to be in my thesis committee.

I would like to specially thank two of the members of our group: Dr. Martin Indarte for patiently analyzing thousands of molecules visually and teaching me the tricks of visual inspection, and Dr. Lionel Costenaro for doing the binding assays for human bleomycin hydrolase.

Many of my colleagues at IBB became friends over the years. I have been lucky to work in a place where people spend time together outside office hours. I would like to thank Oscar, Sasha, Pau, Isaac, Juan, Antonio, Ricard, Manolo, Dunja and especially to Rita for sharing the miseries of doing a PhD.

I also want to thank all my family and friends. My mother Perihan, my father Kamil and my sister Gulin who supported and believed in me no matter what and I am lucky to have them. I want to thank to Gunes, who never left me alone in Barcelona whenever I needed her; she is a true friend. Finally, I am going to thank to myself for bearing with me and teaching me I can do anything I set my mind to.





## ***Abstract***

Virtual screening plays a central role in the world of drug discovery today. In silico testing allows to screen millions of small molecules and to choose only the most promising ones for experimental testing. To find potential drug candidates, it is crucial to bring together individual and complementary computational tools. In this thesis, I describe an automated virtual screening procedure that combines pharmacophore modeling and searches, high-throughput molecular docking, consensus scoring and binding free energy estimation with the linear interaction energy (LIE) method through molecular dynamics simulations.

One goal of this thesis was to build an evolving and versatile virtual screening methodology, which enables integration of different tools at different steps. The procedure that started as a combination of a simple size filter, molecular docking and consensus scoring, advanced into an elaborate and automated computational workflow with the addition of pharmacophore searches and binding free energy estimation with LIE. This integrated method intends to compensate for weaknesses of individual structure-based techniques and allows the evaluation and comparison of the performance and accuracy of these techniques. Another important goal was to apply the computational workflow to target proteins and find hits that could be drug candidates. Experimental testing performed for human acid  $\beta$ -Glucosidase and bleomycin hydrolase indicate that several small molecules selected by the computational workflow display micromolar inhibitory activity. The standard LIE method used in this work was applied to more than ten thousand ligand-protein complexes for three different targets, which is, to our knowledge, the first time application of LIE at such large scale.

## ***Resum***

Actualment, el cribratge virtual juga un paper central en el món del descobriment de fàrmacs. L'anàlisi *in silico* permet el cribratge de milions de molècules petites i la tria de les més prometedores per a les proves experimentals. Per trobar candidats que puguin esdevenir fàrmacs, és crucial reunir una sèrie d'eines computacionals individuals i complementàries. En aquesta tesi, es descriu un procediment automatitzat de cribratge virtual que combina el modelat de farmacòfors i el seu ús en cerques, mètodes d'alt rendiment d'acoblament molecular, puntuació de consens i estimació d'energia lliure d'unió mitjançant el mètode d'energia d'interacció lineal (LIE) a partir de simulacions de dinàmica molecular.

Un dels objectius d'aquesta tesi ha estat el de construir una metodologia flexible i versàtil de cribratge virtual, que permeti la integració de diferents eines en les diferents etapes de l'estudi. El procediment, que es va iniciar com la combinació d'un senzill filtre per tamany, la simulació de l'acoblament molecular i una puntuació de consens, ha derivat en un procediment computacional elaborat i automatitzat amb l'addició de cerques basades en farmacòfor i l'estimació de l'energia lliure d'unió mitjançant el mètode LIE. Aquest mètode integrat té l'objectiu de compensar les debilitats individuals de les diferents tècniques usades i permet avaluar i comparar el rendiment i la l'exactitud d'aquestes tècniques. Una altra fita important ha estat l'aplicació del procediment computacional a proteïnes diana concretes per tal d'avaluar-ne la capacitat de trobar molècules que puguin ser candidats a fàrmacs. Tests experimentals realitzats

per a la  $\beta$ -Glucosidasa àcida i la hidrolasa de Bleomicina humanes indiquen que diverses molècules petites seleccionades pel procediment computacional tenen activitat inhibidora micromolar. El mètode LIE emprat en aquest treball es va aplicar sobre més de deu mil complexos proteïna-ligand per a tres proteïnes diana diferents, el que és, al nostre entendre, la primera aplicació del mètode LIE a aquesta escala.

# Contents

<i>List of Figures</i> .....	xiii
<i>List of Tables</i> .....	xvii
<i>List of Equations</i> .....	xix
<i>List of Abbreviations</i> .....	xxi
<i>Preface</i> .....	1
1. <i>Drug Design Introduction</i> .....	5
1.1. <i>High-Throughput Screening (HTS)</i> .....	5
1.2. <i>Virtual Screening (Computer Aided Rational Design)</i> .....	6
1.3. <i>Ligand-Based Methods</i> .....	7
1.3.1. <i>Ligand-based Pharmacophore Modeling</i> .....	7
1.3.2. <i>QSAR</i> .....	8
1.4. <i>Receptor-Based Methods</i> .....	9
1.5. <i>Configuration Generation Problem</i> .....	9
1.5.1. <i>Receptor-based Pharmacophore Modeling and Searching</i> .....	10
1.5.2. <i>Docking</i> .....	11
1.6. <i>Affinity prediction problem</i> .....	12
1.6.1. <i>Factors determining ligand-receptor binding affinity</i> .....	13
1.7. <i>Approaches for prediction of binding affinities</i> .....	17
1.7.1. <i>Scoring Functions and Consensus Scoring</i> .....	18
1.7.2. <i>FEP and TI</i> .....	20
1.7.3. <i>MM-PBSA</i> .....	21
1.7.4. <i>LIE</i> .....	22
2. <i>Objectives</i> .....	25
3. <i>The Three-Step Docking Procedure versus the Hybrid Procedure</i> .....	29
3.1. <i>Summary</i> .....	29
3.2. <i>Biological Background</i> .....	29
3.2.1. <i>Sample Protein: Human Endothelial Nitric-Oxide Synthase (eNOS)</i> .	29
3.2.2. <i>Benchmark Protein: Human Checkpoint Kinase 1</i> .....	30
3.3. <i>Common Methods For Three-Step Docking And The Hybrid Procedure</i> .	32
3.3.1. <i>Preparation Of The Small Molecule Libraries And The Proteins</i> .....	32
3.3.2. <i>Collection Of The Known Actives</i> .....	33
3.3.3. <i>Consensus Scoring</i> .....	36
3.4. <i>Methods For The Three-Step Docking Approach</i> .....	37
3.4.1. <i>Parameter Calibration For Docking With The Sample Protein</i> .....	37

3.4.2. <i>The Size Filter</i> .....	40
3.4.3. <i>Docking Experiments For The Three-Step Docking Approach</i> .....	40
3.5. <i>Methods for the Hybrid Approach</i> .....	41
3.5.1. <i>Pharmacophore Filtering</i> .....	42
3.5.2. <i>Docking with AutoDock Vina for the Hybrid Method</i> .....	44
3.6. <i>Results And Discussion</i> .....	44
3.7. <i>Conclusion</i> .....	50
4. <i>The Computational Workflow</i> .....	54
4.1. <i>Biological Background</i> .....	54
4.1.1. <i>Human T-Protein</i> .....	54
4.1.2. <i>Human Bleomycin Hydrolase</i> .....	55
4.1.3. <i>Human Acid <math>\beta</math>-Glucosidase</i> .....	56
4.2. <i>Methods</i> .....	57
4.2.1. <i>Preparation Of Small Molecule Libraries And The Target Proteins</i> ..	58
4.2.2. <i>Pharmacophore Creation And Search</i> .....	58
4.2.3. <i>Docking the pharmacophore-filtered libraries to the targets</i> .....	64
4.2.4. <i>Consensus Scoring</i> .....	65
4.2.5. <i>LIE</i> .....	66
4.3. <i>Results and Discussion</i> .....	68
4.3.1. <i>Results for human T-protein and human Bleomycin Hydrolase</i> .....	68
4.3.2. <i>Results for human GCase</i> .....	87
4.4. <i>Conclusion</i> .....	112
4.5. <i>Future Work</i> .....	113
5. <i>Conclusions</i> .....	116

## List of Figures

Figure 1: A representation of the Lennard-Jones 12-6 function.....	15
Figure 2: Examples of non-covalent interactions found in protein-ligand complexes. .	16
Figure 3: The role of water in protein ligand binding .....	17
Figure 4: Structure of eNOSox binding site co-crystallized with 6-nitroindazole (active site inhibitor orientation).....	30
Figure 5: Schematic representation of the cell cycle .....	31
Figure 6: Catalytic residues and binding site of Chk1.....	32
Figure 7: Known actives of Chk1.....	35
Figure 8: The energy differences between group parameter sets and the reference set.	39
Figure 9: The flowchart summarizing the three-step docking algorithm.....	41
Figure 10: The flowchart summarizing the hybrid method.....	42
Figure 11: The pharmacophore filter created from Chk1 binding site residues Cys87 and Asn135. ....	43
Figure 12: The pharmacophore filter created from three known binders of Chk1.....	44
Figure 13: Normalized scores after the first docking step of the three-step docking approach before 0.5 % truncation.....	45
Figure 14: Normalized scores after the first docking step of the three-step docking approach after 0.5 % truncation. ....	45
Figure 15: Docked conformations of the known ligands that passed the first pharmacophore filter.....	49
Figure 16: Docked conformations of the known ligands that passed the second pharmacophore filter.....	50
Figure 17: Glycine cleavage system mechanism.....	55
Figure 18: Cysteine protease mechanism of peptide bond cleavage .....	56
Figure 19: The binding site of human T-protein with the competitive inhibitor bound. .	59
Figure 20: The pharmacophore filter created from the T-protein binding site residue Glu204 and from the known inhibitor .....	60
Figure 21: The binding site of human bleomycin hydrolase .....	61
Figure 22: The pharmacophore filter created from the bleomycin hydrolase binding catalytic site residues Cys73 and His372 .....	61
Figure 23: Binding site of 2NSX and IFG.....	62
Figure 24: The first pharmacophore filter (pharma1) designed for GCCase .....	63
Figure 25: The second pharmacophore filter (pharma2) designed for GCCase.....	64
Figure 26: Normalized scores calculated without truncation with the five scoring functions.....	69
Figure 27: Normalized scores calculated with the second normalization procedure against compound rank (after the 0.5% truncation) for human T-protein.....	70

<i>Figure 28: Normalized scores calculated without truncation with the five scoring functions.....</i>	<i>71</i>
<i>Figure 29: Normalized scores calculated with the second normalization procedure against compound rank (after the 0.5% truncation) for human bleomycin hydrolase... </i>	<i>72</i>
<i>Figure 30: Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for human T-protein.....</i>	<i>74</i>
<i>Figure 31: Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for human bleomycin hydrolase.....</i>	<i>75</i>
<i>Figure 32: Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) for human T-protein. ....</i>	<i>76</i>
<i>Figure 33: Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) for human bleomycin hydrolase. ....</i>	<i>76</i>
<i>Figure 34: Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) versus the ligand size for human T-protein.....</i>	<i>77</i>
<i>Figure 35: Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) versus the ligand size for human bleomycin hydrolase.....</i>	<i>78</i>
<i>Figure 36: The distribution of water molecules within 4 Å of the ligands at the end of docking experiments and LIE simulations for human T-protein.....</i>	<i>79</i>
<i>Figure 37: The distribution of water molecules within 4 Å of the ligands at the end of docking experiments and LIE simulations for human bleomycin hydrolase.....</i>	<i>79</i>
<i>Figure 38: The competitive inhibitor 5-CH3-H4-folate in the binding site (comparison with docking) .....</i>	<i>80</i>
<i>Figure 39: The competitive inhibitor 5-CH3-H4-folate in the binding site (comparison with LIE).....</i>	<i>81</i>
<i>Figure 40: Configurations of the competitive inhibitor 5-CH3-H4-folate superimposed (actual, docking and LIE binding modes). ....</i>	<i>81</i>
<i>Figure 41: The molecules selected for experimental testing of inhibition of human bleomycin hydrolase.....</i>	<i>83</i>
<i>Figure 42: Binding mode of Compound 1 as predicted by docking.....</i>	<i>85</i>
<i>Figure 43: Binding mode of Compound 2 as predicted by docking.....</i>	<i>85</i>
<i>Figure 44: Binding mode of Compound 5 as predicted by docking.....</i>	<i>86</i>
<i>Figure 45: Binding mode of Compound 7 as predicted by docking.....</i>	<i>86</i>
<i>Figure 46: Binding mode of Compound 9 as predicted by docking.....</i>	<i>87</i>
<i>Figure 47: Normalized scores calculated without truncation with the five scoring functions for dock1 experiment of GCase.....</i>	<i>89</i>
<i>Figure 48: Normalized scores calculated with the second normalization procedure against compound rank (after the 0.5% truncation) for dock1 experiment of GCase. ..</i>	<i>90</i>
<i>Figure 49: Normalized scores calculated without truncation with the five scoring functions for dock2 experiment of GCase.....</i>	<i>91</i>

<i>Figure 50: Normalized scores calculated with the second normalization procedure against compound rank (after the 0.5% truncations) for dock2 experiment of GCase.</i>	92
<i>Figure 51: Normalized scores calculated without truncation with the four scoring functions for dock3 experiment of GCase.</i>	93
<i>Figure 52: Normalized scores calculated with the second normalization procedure against compounds rank (after the 0.5% truncation) for dock3 experiment of GCase.</i>	94
<i>Figure 53: Numbers of intersecting molecules chosen by more than one docking experiment in the top 600.</i>	95
<i>Figure 54: Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for dock1.</i>	96
<i>Figure 55: Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for dock2.</i>	97
<i>Figure 56: Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for dock3.</i>	98
<i>Figure 57: Rmsd differences between the docking results and LIE simulation results for GCase.</i>	99
<i>Figure 58: Rmsd differences between the docking ligand configurations and LIE simulation results versus the ligand size for GCase.</i>	100
<i>Figure 59: The distribution of water molecules within 4 Å of the ligands at the end of docking experiments and LIE simulations for Gcase.</i>	101
<i>Figure 60: Molecules chosen for experimental testing.</i>	103
<i>Figure 61: Hydrolysis of the substrate by the imiglucerase activity.</i>	104
<i>Figure 62: Experimentally observed activities of the selected molecules.</i>	105
<i>Figure 63: Docked binding modes of five hits.</i>	107
<i>Figure 64: Docked and simulation binding modes of Compound 3.</i>	108
<i>Figure 65: Docked and simulation binding modes of Compound 4.</i>	109
<i>Figure 66: Docked and simulation binding modes of Compound 13.</i>	109
<i>Figure 67: Docked and simulation binding modes of Compound 19.</i>	110
<i>Figure 68: Docked and simulation binding modes of Compound 21.</i>	111





## List of Tables

<i>Table 1: List of drugs that were discovered by rational design</i> .....	6
<i>Table 2: The references for the molecules chosen for testing the methods</i> .....	34
<i>Table 3: Parameter values used for docking of the test set of molecules to human eNOS</i> . .....	38
<i>Table 4: Molecule groups divided according to the number of rotatable bonds</i> .....	38
<i>Table 5: Relative ranking of known molecules with 0.5% truncation and without any truncation</i> .....	46
<i>Table 6: Ranks of known ligands at the end of the first docking step (after 0.5% truncation)</i> .....	47
<i>Table 7: Ranks of known ligands that passed the first pharmacophore filter (after 0.5% truncation)</i> .....	48
<i>Table 8: Ranks of known ligands that passed the second pharmacophore filter (after 0.5% truncation)</i> .....	48
<i>Table 9: Parameters used for free state simulation of the ligand in a sphere filled with water in the initial equilibration period</i> .....	68
<i>Table 10: Parameters used for bound state simulation of the ligand in the protein binding site in the initial equilibration period</i> .....	68
<i>Table 11: Pearson's Correlations of scoring functions with each other and normalized consensus score of the selected molecules for human T-protein</i> .....	73
<i>Table 12: Pearson's Correlations of scoring functions with each other and normalized consensus score of the selected molecules for human bleomycin hydrolase</i> .....	74
<i>Table 13: The list of selection criteria fulfilled by the chosen molecules for testing against human bleomycin hydrolase</i> .....	82
<i>Table 14: Observed activities of the candidate compounds for human bleomycin hydrolase</i> .....	84
<i>Table 15: Pearson's Correlation coefficients between different scoring functions and NCS99.5 for docking experiment dock1</i> .....	88
<i>Table 16: Pearson's Correlation coefficients between different scoring functions and NCS99.5 for docking experiment dock2</i> .....	90
<i>Table 17: Pearson's Correlation coefficients between different scoring functions and NCS99.5 for docking experiment dock3</i> .....	92
<i>Table 18: Pearson's Correlations of scoring functions with each other and with LIE energies for the selected molecules from dock1 experiment</i> .....	96
<i>Table 19: Pearson's Correlations of scoring functions with each other and with LIE energies for the selected molecules from dock2 experiment</i> .....	97
<i>Table 20: Pearson's Correlations of scoring functions with each other and with LIE energies for the selected molecules from dock3 experiment</i> .....	98
<i>Table 21: The list of selection criteria fulfilled by the chosen molecules</i> .....	104
<i>Table 22: Volume ratio for activity studies of imiglucerase</i> .....	105

<i>Table 23: Experimentally observed activities of the selected molecules.....</i>	<i>107</i>
<i>Table 24: Hydrogen bonds made between the hit compounds and GCase. ....</i>	<i>112</i>

## **List of Equations**

<i>Equation 1: Equation for ligand-receptor binding.....</i>	<i>13</i>
<i>Equation 2: Equation for binding free energy.....</i>	<i>13</i>
<i>Equation 3: Equation for <math>K_i</math>.....</i>	<i>13</i>
<i>Equation 4: Equation for Coulombic interactions.....</i>	<i>14</i>
<i>Equation 5: Equation for Lennard-Jones 12-6 function.....</i>	<i>14</i>
<i>Equation 6: Energy equation for FEP.....</i>	<i>20</i>
<i>Equation 7: Energy equation for TI.....</i>	<i>21</i>
<i>Equation 8: Energy equation for MM-PBSA.....</i>	<i>21</i>
<i>Equation 9: Energy equation for LIE.....</i>	<i>22</i>
<i>Equation 10: Equation for normalized score.....</i>	<i>36</i>
<i>Equation 11: Equation for normalized consensus score.....</i>	<i>36</i>
<i>Equation 12: Equation for calculating the the solvation sphere radius for LIE.....</i>	<i>66</i>



## *List of Abbreviations*

3D QSAR	Three Dimensional Quantitative Structure Activity Relationship
6NI	6-Nitroindazole
AMBER	Assisted Model Building And Energy Refinement
Chk1	Human Checkpoint Kinase
CNS	Central Nervous System
eNOS	Endothelial Nitric-Oxide Synthase
ER	Endoplasmic Reticulum
ERAD	Endoplasmic Reticulum Associated Degradation
ERT	Enzyme Replacement Therapy
FEP	Free Energy Perturbation
GBA	Glucosidase, Beta, Acid Gene
GCase	Acid B-Glucosidase
GCS	Glycine Cleavage System
GD	Gaucher's Disease
GROMOS	Groningen Molecular Simulation Package
HTS	High Throughput Screening
IFG	Isofogamine
iNOS	Inducible Nitric-Oxide Synthase
LIE	Linear Interaction Energy
MD	Molecular Dynamics
MM-PBSA	Molecular Mechanics Poisson Boltzmann Surface Area
NCS	Normalized Consensus Score
NMDA	N-Methyl-D-aspartate
NMR	Nuclear Magnetic Resonance
NN-DNJ	N-nonyl-deoxynojirimycin
nNOS	Neuronal Nitric-Oxide Synthase
NO	Nitric Oxide
OPSL	Optimized Potential For Liquid Simulations
PDB	Protein Databank
QSAR	Quantitative Structure Activity Relationship
RMSD	Root Mean Square Deviation
SMILES	Simplified Molecular Input Line Entry Specification
SRT	Substrate Reduction Therapy
TI	Thermodynamic integration
VSL-1	Virtual Screening Library 1
VSL-2	Virtual Screening Library 2



# *Preface*

## *What is all this about?*

The world of drug design is nothing but a huge maze with maybe thousands of dead-ends trapped with bad strategies, disguised as seemingly “good ideas”, that cannot be predicted beforehand, and a single exit—the existence of which is open to discussion. It is a very complex world that includes a roughly estimated number of 500,000 proteins, only around 10,000 of which have been characterized structurally, generated by an estimated 20,000 to 25,000 open reading frames in the human genome.<sup>1</sup> Also, the project of developing a new drug may need tens of researchers working for up to ten years only to reach to the point of clinical trials, costing around half a billion dollars and may still fail miserably for many different reasons.<sup>2</sup> So, in this world of unknowns, estimations and an overwhelming probability of failure, is there a magical way that would solve the problem of drug design and if there is, are we close to find it? This is a question that cannot be answered till this magical way is found and till then, what can be done is, basically, trial and error.

With the hope and aim of surmounting some of the obstacles found at the early stages of drug discovery, an automated computational workflow was created in this study. In the following introductory chapter, I begin by introducing the most important concepts and methodologies of computational drug design. The second chapter summarizes my motivations and goals for pursuing computational drug design research. The third chapter tells the story of a simple computational methodology emerging and evolving into a more complicated approach combining widely used computational drug design techniques. The fourth chapter is about the establishment and application of the final and automated version of the methodology explained in chapter three, with the addition of molecular dynamics simulation and binding free energy prediction by the linear interaction energy method. Finally, the last chapter summarizes and concludes the overall study. A successful computational drug design procedure needs several methodologies combined, and the workflow developed and evolved in this thesis aims to become a useful tool by doing so.





# *Chapter 1*

## *Drug Design Introduction*



# ***1. Drug Design Introduction***

Biological systems are the most complex type of systems we know and molecular recognition between small molecules (ligands) and their target proteins (receptors) lie at the center of many of the processes taking part in biological systems. The rapid increase in the number of receptor proteins with known three-dimensional structure has opened the possibility to discover the accurate binding poses of their natural or pharmacological ligands in their binding sites, one of the keys towards understanding the proteins' function mechanisms. Protein activity is often modulated by binding of a small molecule to specific sites on the protein, and the disturbance of such regulation is mostly a cause for a disease.<sup>2</sup> Hence, many drugs work by inhibiting the function of proteins with enhanced or unbalanced activity. Also, the proper activity of malfunctioning or non-functioning proteins causing a disease can sometimes be restored by the binding a specific chemical.

The development of a new drug is a highly expensive, difficult and painstaking process, which can last for years and consists of several steps that start with a search for a candidate ligand with a noticeable affinity for the target protein in question. This step is called “lead discovery” and outputs a ligand to be further optimized for an increased affinity and selectivity. With increased affinity and selectivity, the candidate drug should also show optimum pharmacokinetic properties, including its absorption, distribution and metabolism in the body, along with its excretion and lack of toxicity. It is after all these criteria are met that the drug enters several steps of clinical testing before being validated and marketed.

While finding a potent and high-affinity lead ligand in a fast and reliable way is already a challenge quite hard to tackle with as is, the significant increase in the number of therapeutic targets without known small molecule ligands made available in the current “post-genome era” has made the search for a lead structure even a bigger challenge.<sup>3</sup> In the past, most drugs were discovered either by identification of the active ingredient from traditional remedies, modification of natural ligands or by serendipitous discovery. However, new discovery approaches are based on understanding the molecular and physiological control mechanisms of the disease. In the big quest for a small lead molecule, there are two major approaches to ease these entanglements: high-throughput (experimental) screening<sup>4,5</sup> –*in vitro* testing– and virtual screening (rational design)<sup>6,7</sup> –*in silico* testing– of large compound libraries.

## ***1.1. High-Throughput Screening (HTS)***

High-throughput screening (HTS) is an experimental random screening method in drug discovery and it involves testing of large molecule libraries composed of natural or synthetic compounds for possible biological activity, independent of their actual chemical properties. The use of robotics systems has boosted the capability to conduct millions of chemical or pharmacological tests, enabling synthesis of thousands of compounds in a short time from a few reagents and rapid identification of active compounds. Even though HTS methodologies can be considered being unbiased without preconceived restrictions about the tested compounds, ignoring the biological features of the target makes them “irrational” due to the random and untargeted screening of molecules for as many as possible and as fast as possible.

Consequently, for more efficient experimental screening, targeted approaches that involve preselection of compounds by computer methods based on “drug-likeness” have been developed.<sup>8</sup> By this way, it has become possible to promptly identify and eliminate candidate molecules that are unlikely to survive the later stages of discovery and development.

### 1.2. Virtual Screening (Computer Aided Rational Design)

Virtual screening is a computational method to scan large numbers of small molecules to see whether they bind to a target protein and function in the desired manner, using the available information about the target, binding mechanism or the known or hypothetical binding mode.<sup>7</sup> It is a complementary approach to experimental discovery methods that aims to enhance and accelerate the lead discovery process. Drug discovery research for both hit identification and lead optimization has shifted towards computational methodologies, which are able to handle millions of molecules in a much shorter time compared to experimental techniques/approaches. The increase in the number of known protein structures and the enormous chemical space of conceivable small molecules has drawn particular attention to virtual screening techniques.<sup>6</sup>

Even though virtual screening is a newly emerging approach, the advances in computer technology and methodology promoted its success, and there are already several drugs that were developed and optimized fully or partially with rational design techniques (Table 1).

Drug	Target	Disease or Infection
Dorzolamide <sup>9</sup>	Carbonic anhydrase	Glaucoma
Imatinib (Gleevec) <sup>10,11</sup>	Tyrosine kinase	Some types of cancer
Cimetidine <sup>12</sup>	Histamine H2 receptor	Peptidic ulcers
Zolpidem <sup>13</sup>	GABAA receptor	Insomnia, brain disorders
Zanamivir <sup>14</sup>	Neuraminidase	Prophylaxis of influenza
Raltegravir <sup>15</sup>	HIV integrase	HIV infection
Enfuvirtide <sup>16</sup>	HIV transmembrane protein	HIV infection

**Table 1:** List of drugs that were discovered by rational design, or where rational design played a key role in the discovery process.

Drug discovery is such a difficult problem that every relevant technique has to be utilized to its best advantage. All computational techniques may provide different strategies, useful insights, new suggestions for molecular structures to synthesize, and cost-effective virtual analysis prior to synthesis.<sup>1</sup> The strategy to be pursued in rational design strictly depends on the availability of the three-dimensional structure of the biological target. Therefore, computer-aided drug discovery techniques can be grouped in two classes: ligand-based and receptor-based (often also called target-based) methods.<sup>17,18</sup>

### ***1.3. Ligand-Based Methods***

Despite the increase in the number of proteins with known 3-D structure, there are still a fair number of drug targets without the structure information. However, if there is at least one known active ligand validated through a cell culture assay for a target, ligand-based computational techniques that do not require the target structure and binding site geometry can be employed. The main idea followed in ligand-based drug design is that ligand structural similarity or similarity of steric and electrostatic features implies similar activity, which allows deriving the required properties for active molecules from the analysis of already known active ligand(s).<sup>19</sup> While not requiring a target structure can be advantageous for ligand based methods, with respect to receptor-based methods, it can at the same time be a drawback not to be able to integrate ligand-target complementarity information in the drug design process. On one hand ligand based methods have low computational complexity; however, on the other hand they do not allow a chemically diverse set of results due to restriction by the known ligand(s). The ligand based techniques range from rather simple similarity searches<sup>20</sup> –usually applied if there is only one known active ligand– to more sophisticated methods like pharmacophore modeling<sup>21</sup> or statistical methods (QSAR)<sup>22</sup> in cases where several active compounds are known.

#### ***1.3.1. Ligand-Based Pharmacophore Modeling***

In computer-aided drug design, one approach to distinguish potentially active from inactive compounds in a database of small molecules is to use the knowledge of the physical and chemical properties of the target binding site or a set of known actives.<sup>21,23</sup> Pharmacophores are ensembles of these physical and chemical features that are necessary for optimal interactions between a specific biological target and a ligand to enhance or inhibit the target function. The most common pharmacophoric features include being aromatic, hydrophobic, hydrogen bond donor, hydrogen bond acceptor, an anion or a cation. Therefore, based on which pharmacophoric features are used, pharmacophore modeling approaches can be categorized in two, whether the structural properties of the target protein and/or the binding site are known, namely structure-based pharmacophore models, or a set of known active ligands which bind to the same region in the target protein are known, namely ligand-based pharmacophore modeling. In this section, we concentrate on the ligand-based approach; the receptor-based pharmacophore models are discussed in section 1.5.1.

For ligand-based pharmacophore modeling, the information of the target protein or the binding site is not needed; the model can be created from a set of known actives. However, it is crucial that all the known ligands should bind to the same region of the target protein. Ligand-based pharmacophore model creation is basically finding the common chemical and physical features of the known ligands to be used as a “query” to search for molecules fitting the model in a small molecule database.<sup>24</sup>

If there are no data available about binding conformations of the known ligands to the target protein, finding the active conformations may be quite challenging.<sup>1</sup> In this case, all the conformers of the ligands should be created and aligned to find the best alignment. The alignment can be done first by superimposing the most rigid compounds and then adjusting the remaining compounds accordingly with computational tools or manually. To prevent any

bias, the known molecules should preferably be derivatives of different structures. Once the alignment is completed, the pharmacophore model can be generated via determining the features that are present in all the molecules. The model can then be used to search a database of molecules, resulting in a qualitative ranking based on how well the molecules fit the model.

### 1.3.2. QSAR

QSAR stands for “quantitative structure-activity relationship” and this technique can be categorized into the conventional QSAR and the 3-D QSAR according to the process followed, even though the concept is similar for both.<sup>25</sup> The quantitative structure-activity relationships are basically a set of correlations between the molecular composition and structure of a compound and its biological or chemical activity, and the derivation of a simple equation that combines these correlations. The equation coefficients, which are weights of molecular properties, are derived by curve fitting. These molecular properties are called “descriptors” and they can be any numerical value that describes the molecule. In classical QSAR, the descriptors can be structural features or physicochemical or steric properties such as molecular weight, number of hydrogen bond donor or acceptor atoms, number of heavy atoms, number of rotatable bonds etc.<sup>25</sup>

3D QSAR methods, on the other hand, analyze three-dimensional structures and binding modes and affinity of the ligands to an active site in one specific target. 3D QSAR methods attempts to define the properties of an active site, without actually knowing its structure, through the computation of steric and electrostatic interactions between a known active ligand and putative probe atoms placed at various positions on a grid surrounding the known ligand.<sup>25</sup>

The classical QSAR methods predict activity from an equation fit to the descriptors of the known active ligands and their coefficients. Therefore, the first step is to define a training set of known molecules with their experimental activity. The molecules in the training set should be diverse enough to span all the possible values for the activity and also abundant enough to prevent over-fitting to an outlier. After the selection of the molecules for the training set, the descriptors are calculated. Since QSAR models are mostly linear, the correlation coefficients of each descriptor with the activity are calculated to choose which descriptors to include in the QSAR model. The descriptor with the highest correlation coefficient with the activity can be selected. The next descriptor to be selected should have a high correlation with the activity but not a strong correlation with the previously chosen descriptor to prevent redundancy and also to compensate for the weaknesses of the previous descriptor. At this point using a correlation matrix showing the correlations between the descriptors and the activity is quite practical and can help choosing the descriptors. Once the descriptors to be included in the QSAR model are selected, the coefficients of the linear fitting equation are generated. To validate the method, the experimental values can be compared with the values predicted for the training set and the model can be further improved if needed. Thereafter, selected descriptors and generated coefficients can be used to predict the activity of the ligands in the test set.

In 3D QSAR, like the classical QSAR, the first step is to create a training set of known molecules with experimental activity values. However, while 10 known actives are enough for a reasonable model in QSAR, 15-20 known actives are required for an ideal 3D QSAR model. The goal is to find a bioactive conformer that corresponds to the fitting equation of the

classical QSAR. For that, alignments and shapes of the training set molecules are created, paying extra attention to put the molecules in conformations with the highest shape similarity and the closest alignment of the pharmacophoric features. The alignment is done by primarily aligning most rigid molecules and then adjusting the remaining molecules to give the optimum fit to shape and pharmacophoric features.

Overall, any QSAR method can be used if there is a sufficient number of ligands with known experimental activity available. This is one of the main drawbacks of QSAR methods and restricts its applicability. Also, the QSAR models can only successfully predict the activity based on the known active set, meaning they might not be fit for prediction of molecules that are not structurally similar.

#### ***1.4. Receptor-Based Methods***

Receptor-based methods (also called target- or structure-based methods) use structural information about the target, e.g., crystal structures, structures derived by NMR, or homology models.<sup>18</sup> The main assumption of receptor-based design is that good inhibitors must possess significant structural and chemical complementarity to their target receptor.<sup>26</sup> In cases where the information about the target protein structure and/or the binding site is available, receptor-based virtual screening techniques can be applied without any information about known active ligands. Receptor-based virtual screening has become popular due to the increasing number of three-dimensional structures of proteins that may be potential drug targets. In those cases where the target binding site is unknown, binding pocket prediction methods like PASS<sup>27</sup>, PocketPicker<sup>28</sup> or LIGSITE<sup>29</sup> can be employed to predict potential binding sites.

Receptor-based methods have the advantage over ligand-based methods that they use information from the target and provide insight about the mechanism of action. However, they are computationally more expensive and complex because of the existence of the target structure in the system. Proteins are complicated systems for which it is quite difficult to find optimum solvation and force field parameters. The majority of receptor-based methods also assume that the target is rigid or doesn't show significant conformational changes upon ligand binding, thus causing a larger number of false negatives than ligand-based methods due to restrictions on the favorable poses.<sup>30</sup>

The success of a receptor-based drug discovery project relies on two aspects: the generation of reasonable ligand binding modes that would span the entire surface available for the ligand (*configuration-generation problem*) and accurate recognition of the binding modes that would be closest to the experimental situation, and reasonable estimation and ranking of binding affinities of the ligands to the target (*affinity prediction problem*).<sup>2</sup>

#### ***1.5. Configuration Generation Problem***

Ligand-protein interactions lie at the center of most biological processes. It is crucial in computational drug design to accurately estimate these interactions. For a reasonable estimation, the conformational space available for a ligand in the binding site of a protein target should be very well spanned by the ligands tested and the most reasonable conformers should be chosen. There are two prominent ways to explore the conformational space:

pharmacophore searches and docking. These methods can either be used separately or combined in hybrid methodologies.

### ***1.5.1. Receptor-Based Pharmacophore Modeling And Searching***

The alternative to generating pharmacophore models from known ligands is to generate them from the target binding site. This method is preferred if there are not any or enough known ligands and also prevents any errors that might arise from the restrictions of a training set of known ligands. In receptor-based pharmacophore modeling, it is crucial to examine the binding site intensively and deduce the interactions that play major roles in ligand binding and action mechanism of the target. Since the pharmacophoric features are derived from the residues of the active site and the compounds whose pharmacophoric features match the properties of the target binding site are considered to be more active than the other compounds, the matching ligands should have the corresponding pharmacophoric features. For instance, if a hydrogen bond donor feature is defined on a residue of the target binding site, then a corresponding hydrogen bond acceptor feature should be included in the pharmacophore model that is used for small molecule database search.

Even though pharmacophoric features can be automatically derived from the binding site with different tools like Catalyst from Accelrys<sup>31</sup>, MOE from Chemical Computing Group<sup>32</sup>, Phase from Schrödinger<sup>33</sup>, UNITY from Tripos<sup>34</sup>, and LigandScout from Inte:Ligand<sup>35</sup>; visual inspection and manual modification are still advisable. Computational tools may output many possible interactions from the active site residues, but selecting the strongest and most relevant ones is key to a pharmacophore model that would minimize the number of false positives.

One advantage of the receptor-based pharmacophore models to its ligand-based counterpart is that it is possible to define excluded volumes. Excluded volumes are the forbidden regions on the target that cannot be occupied by the ligands. If a ligand causes clashes with the target, then it would be scored with a penalty due to the existence of excluded volume features.

Once the pharmacophore model has been created, the next step is small molecule database search for finding ligands that fit the model. Based on the exhaustiveness of the conformer creation for ligands to be searched, the efficiency and time consumption of the pharmacophore search can change dramatically. Only ground state conformations of the molecules can be stored in a database; however, there is no guarantee that these are the biologically active conformations and they are checked against the pharmacophore model.<sup>1</sup>

The most plausible solution for the *configuration-generation problem* is to search all possible conformations of all the ligands and to select the conformations that fit best to the pharmacophore model. This can be done either by generating all possible conformers and creating a very large database beforehand and doing the search against this database, or alternatively by storing only one conformer in the database and generating the other conformations as each molecule is searched and finishing the search if a matching conformer is found without searching all possible conformers. The second alternative is both rigorous and practical because it does not only make the search on a single conformer and also finishes the search of a molecule as soon as a matching conformer is found. However, no matter which strategy is employed, multiple conformer searches are still time consuming.

Despite all the potential fallbacks and restrictions, pharmacophore models and searches are still valuable techniques in computational drug discovery, especially when used in



combination with high-throughput docking, either as filters to narrow down the number of molecules to be docked or as evaluation methods for the docked ligands in the target binding site.<sup>36</sup>

### 1.5.2. Docking

Docking, which involves a simulation of binding of all molecules in a database to the actual or potential binding site of a target protein, may be the most prominently used tool in computational drug discovery studies.<sup>37,38</sup> Since docking calculations simulate the interactions between a ligand and a protein's binding site, and assign a qualitative score to these interactions, the results may be compared to those of biochemical assays.

Docking studies, in general, have two main aims: accurate structural modeling and correct binding affinity predictions. Therefore, most docking algorithms consist of two parts: a search algorithm and a scoring function.<sup>37,39</sup> The search algorithms focus on ligand and sometimes protein flexibility, and explore the conformational space available to the ligand in the binding site. According to how ligand flexibility is treated, search algorithms can be mainly categorized in three: random or stochastic methods, systematic methods and simulation methods.<sup>1</sup> The search algorithm of a docking program can apply only a single method or a combination of different methods.

*Random or stochastic methods* employ random changes on a single ligand or a group of ligands using Monte Carlo or genetic algorithm approaches. Both approaches share similar principles, but differ in application. They both start with an initial single conformation or a set of conformations and proceed by making random changes to the initial set, finally evaluating the newly generated set with a probability function predefined within the algorithm.<sup>37</sup> In Monte Carlo algorithm, first an initial random conformation of a ligand is generated in the target active site. Then this conformation is scored with a scoring function within the algorithm. Afterwards, a random change is made on the ligand conformation and the resulting conformation is scored again. At this point, a Metropolis criterion is used to decide whether the newly obtained conformation is accepted or not. The Metropolis criterion directly accepts the configuration if it scores better than the previous one. When the new configuration does not have a better score, its acceptance probability is given by the Boltzmann factor for the score difference. This search goes on until the desired number of favorable conformations is met. Genetic algorithms, on the other hand, adopt the principles of biological survival and competition and apply the rules of natural evolution to generate solutions for the search algorithm in docking.<sup>37</sup> First, a set of random conformations is generated to form an initial population. This set is scored with a fitness function and the fittest conformations are selected for the production of the next generation. The next generation of conformations is generated through genetic operators: crossover and mutation. When the population size reaches to its limit, each offspring conformation is scored again and the fittest ones are chosen for the following cycle of the algorithm. These steps are repeated until the fixed number of generations is reached or the best fitness score doesn't improve anymore or a solution is found. AutoDock<sup>40</sup> and GOLD<sup>41</sup> implement genetic algorithms.

*Systematic search methods* aim to explore all the degrees of freedom in a molecule without falling into a combinatorial explosion problem.<sup>42</sup> For a systematic conformational search, the

number of possible conformations is directly proportional to the number of rotatable bonds in the molecule and inversely proportional to the size of the rotational angle increment. For a ligand with  $N$  rotatable bonds, if an angle increment of  $\theta$  is used then the number of allowed states will be  $2\pi/\theta$  for each bond, ending up in  $(2\pi/\theta)^N$  possible conformations for a single ligand.<sup>42</sup> To overcome this combinatorial problem, the ligands are incrementally grown into the active site.<sup>37</sup> One approach for incremental search starts with decomposing the ligand into its fragments and docking these fragments separately into the active site followed by linking the best scoring fragments with appropriate linkers.<sup>43</sup> This strategy is also widely used for *de novo* drug design.<sup>44</sup> An alternative approach is to decompose the ligand in rigid and flexible parts followed by docking only the rigid part into the active site and adding the flexible parts incrementally.<sup>45</sup> DOCK<sup>46</sup>, FlexX<sup>47</sup> and Glide<sup>48</sup> use incremental small molecule construction as their search algorithm.

For *simulation methods*, the most popular approach is molecular dynamics; however computation time restrictions are the biggest obstacle.<sup>39</sup> It cannot be guaranteed that the ligands would cross the high energy barriers and leave local minima within the computationally-feasible short simulation time. DOCK, Glide and AutoDock also employ additional simulation methods as complement to their main search algorithms.

AutoDock Vina<sup>49</sup> uses a simulation based method in which the ligand conformational search is started with a random conformation of the ligand, followed by an investigation of the binding site, defined by a grid, by modifying the ligand's coordinates, thus allowing flexibility only on the ligand (although in principle it is also possible to employ flexibility on the amino acid side chains of the binding site of the receptor protein). AutoDock Vina is freely available and it has been shown to give quite reliable results in both predicting binding modes in X-ray structures of protein-ligand complexes and virtual screening.<sup>50,51,52</sup> However, for virtual screening studies, the accuracy of docking is highly dependent on the target and estimating which docking tool is most suitable for a specific target is still not possible.<sup>53,54</sup>

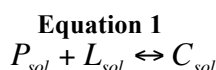
## ***1.6. Affinity Prediction Problem***

The knowledge of 3D structures of the targets and/or the ligands and the accurate prediction of the binding modes and the target-ligand complex structures constitute the basis for receptor-based drug design, however understanding protein-ligand interactions on the molecular level and the accurate prediction of these interactions determine the true success of a receptor-based drug design project. Even though different docking algorithms are able to produce experimentally observed binding modes of ligands to a protein, it is still a challenge to recognize and pick them in huge libraries and assign accurate scores to rank them.<sup>55</sup> The efficiency of a computational drug design procedure relies on accurate prediction of binding affinities.<sup>2</sup>

### 1.6.1. Factors Determining Ligand-Receptor Binding Affinity

Ligands bind either covalently or non-covalently to their targets, based on their structural and energetic recognition with the target.<sup>2</sup> However, the majority of the currently available drugs act via non-covalent interactions with their target proteins.<sup>24</sup> This makes non-bonded interactions of particular interest in drug design and triggers a special interest for finding computational methods to predict ligand-target interactions at the atomic level and the binding affinity of the ligand.

The non-covalent and reversible binding of a ligand (L) to a receptor protein (P) to form a complex (C) almost always takes place in a solution (Equation 1).



The binding affinity is determined by the Gibbs free energy of binding ( $\Delta G$ ) and is related to the experimentally measured binding constant,  $K_i$  (Equation 2 and Equation 3), where  $R$  is the ideal gas constant and  $T$  the temperature.  $\Delta G$  is composed of two parts: an enthalpic ( $\Delta H$ ) and an entropic ( $T\Delta S$ ) component. In,  $[C]$ ,  $[R]$  and  $[L]$  represent the molar concentrations of the complex, protein and the ligand respectively.

The experimentally determined range of the binding constant,  $K_i$  is between  $10^{-2}$  and  $10^{-12}$  M, which corresponds to a Gibbs free energy of binding,  $\Delta G$  between  $-10$  and  $-70$  kJ/mol in solution at  $T=298$  K.<sup>56</sup>

$$\text{Equation 2}$$
$$\Delta G = -RT \ln K_i = \Delta H - T\Delta S$$

$$\text{Equation 3}$$
$$K_i = \frac{[P][L]}{[C]}$$

The binding affinity is used to describe how strongly a ligand binds to its target and is dominated by non-covalent interactions such as electrostatic and van der Waals forces, including solvation and desolvation contributions.<sup>24</sup> These interactions are crucial for structural and energetic recognition between a ligand and a target. Although non-covalent interactions are way weaker than covalent bonds, small stabilizing interactions accumulate to make important contributions to stabilize ligand binding. For a tight ligand binding to a target, some requirements should be fulfilled:<sup>24</sup>

- There should be a high level of steric complementarity between the ligand and the target protein.
- The surface properties of the ligand and the target protein should chemically complement each other. Since lipophilic parts of the protein are mostly in contact with the lipophilic parts of the ligand and polar groups are usually paired accordingly to form hydrogen bonds or ionic interactions, surface properties of both sides should match.

- The ligand should be in an energetically favorable conformation for stability.

### *Electrostatic interactions*

Electrostatic interactions give rise to the forces that are generally accepted to be the most important intermolecular driving factors for ligand-protein binding.<sup>2</sup> Electrostatic phenomena in biomolecular systems are very complex due to the long-range nature of electrostatic forces between large numbers of interacting atoms. The presence of charged groups in proteins further complicate the scenario.

The electrostatic potential energy is represented as a pair-wise summation of Coulombic interactions (Equation 4).

**Equation 4**

$$E_{coul} = \sum_{i=1}^N \sum_{j>i}^N \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}}$$

In Equation 4,  $N$  is the number of atoms in the system,  $q$  is the charge on each atom,  $r$  is the distance between each atom pair  $i$  and  $j$  and  $\epsilon_0$  is the dielectric constant of the environment. The double sum in this interaction function is generally simplified, with a variety of approximations that basically reduce the number of pairs to a tractable one.

### *Van der Waals Interactions*

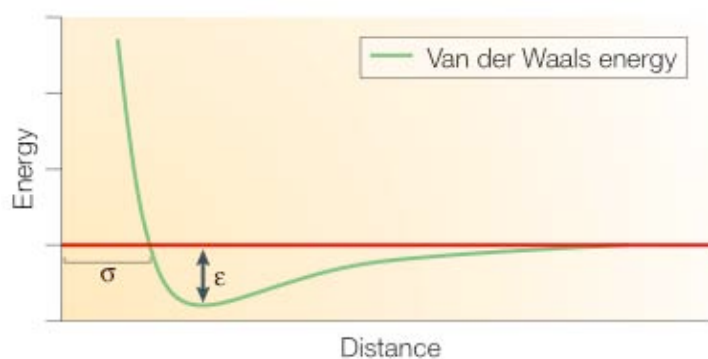
van der Waals interactions have an attractive and a repulsive component.<sup>57</sup> The fluctuations of the electronic charge distribution around the atoms arising from the correlated movements of electrons in interacting molecules (instantaneous polarisation) cause the van der Waals attraction or dispersion force. Repulsion dominates the interaction at short distances and is due to the exclusion principle that prevents the overlap of electron orbitals. As two molecules come closer, the attraction increases until they are separated by the van der Waals *contact distance* (minimum of the potential energy, see Figure 1). Below this distance, repulsion quickly takes over. van der Waals interactions are short ranged: very weak (attractive) at long distance (above  $3\sigma$ ) and very strong (repulsive) at short distance (below  $\sigma$ ).

The van der Waals energy for non-bonded interactions is often modeled by a Lennard-Jones 12-6 function as shown in Equation 5:

**Equation 5**

$$E_{vdW} = \sum_{j=1}^N \sum_{i=1}^N 4\epsilon \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right]$$

where  $N$  is the number of atoms,  $r_{ij}$  is the distance between atoms  $i$  and  $j$ ,  $\epsilon$  is the well depth of the potential (Figure 1) and  $\sigma_{ij}$  is the distance at which the interaction of atoms  $i$  and  $j$  is zero.



**Figure 1:** A representation of the Lennard-Jones 12-6 function. Small-distance repulsion forces are provided by the  $\exp(12)$  term of the Equation 5, while  $\exp(6)$  term is responsible for the attractive force, which approaches zero as the distance increases. (Figure taken from Kitchen et.al.<sup>37</sup>)

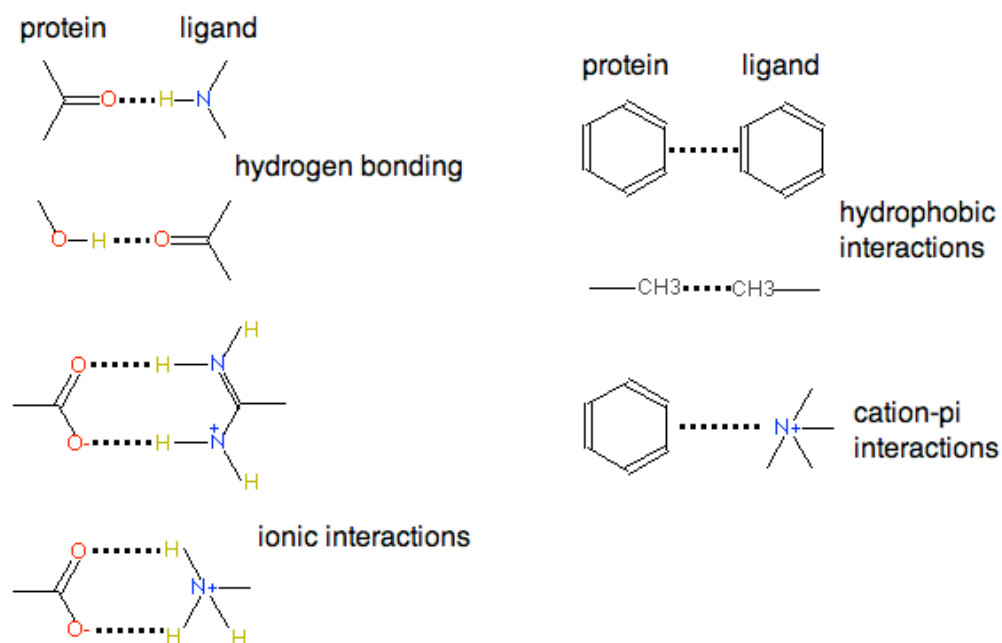
The non-covalent interactions that govern ligand-protein binding are combinations of electrostatic and van der Waals forces and consist of hydrogen bonds (which has been argued to partially have bond character), ionic interactions, hydrophobic interactions, salt bridges,  $\pi$ - $\pi$  interactions and cation- $\pi$  interactions, along with solvation and desolvation (Figure 2).

Hydrogen bond is a form of association resulting from the attraction between an electronegative atom and a hydrogen atom bound to another electronegative atom.<sup>57</sup> The electronegative atoms are mostly, but not exclusively, oxygen, nitrogen or fluorine. Hydrogen bonds between a hydrogen bond donor X—H and a hydrogen bond acceptor Y are formed within distances of 2.5-3.2 Å and angles of 130°-180°.<sup>2</sup> The strength of a hydrogen bond depends highly on its environment, however it is estimated to be less than 20-25 kJ/mol (5-6 kcal/mol) unless a fluorine atom is involved. Hydrogen bonds are weaker than covalent bonds; however, they influence ligand binding strongly by their directional nature.

Salt bridges are ionic interactions formed when oppositely charged moieties (e.g., from a residue side chain and a small ligand) are in close proximity to form an ion pair. They are also capable of making hydrogen bonds; hence they help increase the stability of binding.

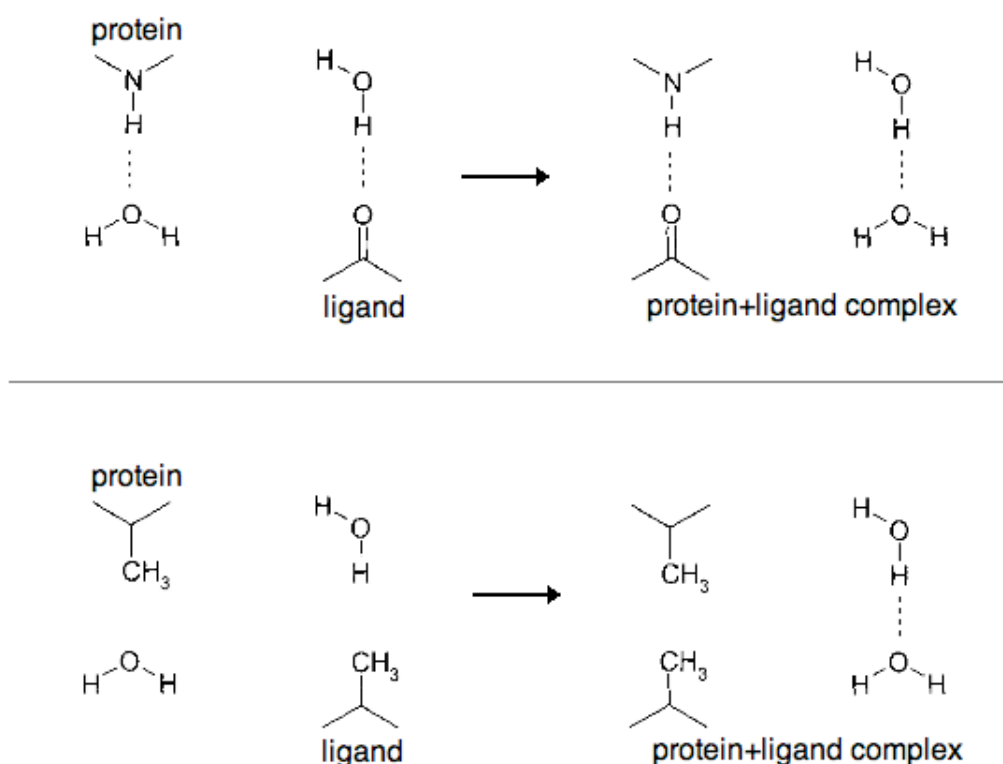
Another significant contribution to ligand-protein binding originates from so-called  $\pi$ - $\pi$  interactions between side chains of residues like tryptophan, phenylalanine or tyrosine and aromatic groups of the ligand.<sup>2</sup>  $\pi$ - $\pi$  interactions lead to stacked arrangements of the aromatic moieties involved.  $\pi$  systems may also interact with cations, an interaction that is observed with relative frequency in ligand-protein binding.<sup>59</sup> The strength of a  $\pi$ -cation interaction is of the same order of magnitude as a hydrogen bond, however it is influenced by different factors, mainly the nature of the cation and the substituents of the  $\pi$  system.<sup>59,60</sup> The exact nature of  $\pi$ - $\pi$  and cation- $\pi$  interactions is still a matter of debate.

Hydrophobic interactions are a consequence of the positive free energy of solvation of apolar groups by water.<sup>58</sup> In the case of protein folding, hydrophobic amino acids such as alanine, valine, leucine, isoleucine, phenylalanine, tryptophan and methionine tend to cluster in a hydrophobic core within the protein that contributes to the stabilization of the folded state. In protein-ligand binding, the same rule applies; hydrophobic moieties of the ligand and the protein try to be in contact, contributing to the binding strength.



**Figure 2:** Examples of non-covalent interactions found in protein-ligand complexes.

Solvation and desolvation make fundamental contributions to both the entropy and enthalpy changes in the system upon protein-ligand binding.<sup>57</sup> Since all biological mechanisms take place in an aqueous environment and ligand binding is no exception, it is highly influenced by solvation and desolvation effects. Therefore the existence of water molecules should be accurately described in protein-ligand binding. In bulk state, each water molecule can make up to four hydrogen bonds, and thus water can strongly influence hydrogen bonding and hydrophobic interactions between the protein and the ligand.<sup>61</sup> In the unbound state, both the ligand and the protein form hydrogen bonds with the water molecules in the environment, and upon complex formation, these water molecules are replaced to allow hydrogen bonding between the protein and the ligand (top part of Figure 3).<sup>24</sup> In the case of hydrophobic interactions (lower part of Figure 3), water molecules are released from the unfavorable environment created by the hydrophobic ligand and the protein residues to form hydrogen bonds.<sup>24</sup> Consequently, hydrophobic interactions between the apolar parts of the ligand and the protein are thermodynamically favoured by the replacement and release of water molecules.



**Figure 3:** The role of water in protein ligand binding. Top part shows how water intervenes in hydrogen bonding while lower part displays the role of water in hydrophobic interactions.<sup>24</sup>

### 1.7. Approaches For Prediction Of Binding Affinities

The success of computational drug design depends on accurate predictions of the binding free energies of small molecules to the target protein.<sup>2</sup> The studies on binding affinity prediction can be categorized in two major groups based on the presence of the knowledge of the 3D structure of the receptor:

- If the 3D structure of the receptor is not known (i.e. ligand-based drug discovery methods), the prediction of the binding affinity of the new ligands is based on the comparison with the known reference ligands with experimentally observed binding affinities. The main assumption for ligand-based affinity prediction is that chemical similarity of the ligands reflects the biological activity.<sup>2,62</sup> One such approach is to compare molecules by considering the presence or absence of functional groups at the one or two-dimensional level, called fingerprinting.<sup>63</sup> Other methods employing topological similarity include substructure mapping<sup>64</sup>, pharmacophore searches<sup>21</sup> and ligand superpositioning<sup>65</sup>. However, these methods only give qualitative binding affinity values. On the other hand, quantitative predictions can be accomplished with the use of QSAR methods.<sup>25</sup> QSAR predicts the binding affinity by finding the correlation between ligands with respect to physicochemical and structural parameters. However two-dimensional QSAR methods suffer from the lack of spatial structure of the individual ligands and lack of receptor interactions. 3D-QSAR, on the other hand,

predicts the binding affinity by correlating the spatial structure of the ligands and experimentally measured binding affinities of the known ligands. Even though 3D-QSAR methods yield reasonable and sensible binding affinity predictions, their dependence on a diverse set of known ligands is still a weak-point.

- If the 3D structure of the receptor is known (receptor-based drug discovery methods), the binding affinity prediction is done based on chemical and geometrical complementarity between the ligands and the target protein. A large variety of methods for binding affinity prediction in receptor-based design have been developed.<sup>2,66</sup> These range from theoretically rigorous to slightly approximate that compromise accuracy for computational speed. Being the fastest, scoring functions are at one extreme of this range, whereas rigorous force-field based methods like free energy perturbation<sup>67</sup> (FEP) calculations or thermodynamic integration<sup>68</sup> (TI) with explicit solvent and flexible ligand and receptor reside at the other extreme. Between these two extremes, but closer to the rigorous side, linear interaction energy<sup>69,70</sup> (LIE) and molecular mechanics Poisson-Boltzmann surface area<sup>71,72</sup> (MM-PBSA) methods are located, both of which have gained considerable attention in recent years.<sup>66</sup>

In this thesis, the focus is on methods used when the 3D structure of the receptor is known.

### ***1.7.1. Scoring Functions And Consensus Scoring***

Scoring functions are approximation methods used to evaluate a docked pose of a ligand to a receptor. Even the most accurate prediction of binding modes of ligands to the target binding site is of little use without a scoring function that produces an accurate ranking of the binding affinities of the test ligands.<sup>55,73</sup> In other words, even if the correct binding conformation of a high-affinity ligand is found by the search algorithm, the result is irrelevant if this compound is not ranked high enough to be selected as a candidate. As for the conformational space search algorithms, a variety of methods are commonly used for scoring in molecular docking and they can be roughly categorized as empirical, force-field based and knowledge based approaches.<sup>37,38</sup> While some scoring functions employ a single approach, some scoring functions use combinations of different approaches.

#### *Empirical scoring functions*

Empirical scoring functions employ fitting to experimental data and define the binding free energy as a sum of parameterized functions obtained from fitting, as first proposed by Bohm.<sup>74</sup> The basic idea behind empirical scoring functions is that the sum of uncorrelated individual terms can be used to approximate the binding affinity. A training set of structurally resolved ligand-protein complexes together with their experimental binding affinities is used in regression analysis to obtain coefficients for individual terms.<sup>37</sup> ChemScore<sup>75</sup> and LUDI<sup>76,77</sup> are scoring functions that employ empirical methods, however individual terms can be handled differently by different scoring functions. For instance, while the hydrogen bonding term in LUDI differentiates between neutral and ionic hydrogen bonds, ChemScore's hydrogen bonding term doesn't.<sup>37</sup> Another point where these two scoring functions differ is the evaluation of the hydrophobic contributions: LUDI uses molecular surface area



representation for the calculation of the hydrophobic term whereas ChemScore just evaluates contacting hydrophobic atom pairs.

Empirical scoring functions are appealing because the terms are usually simple to evaluate, however they can be disadvantageous due to their strong dependence on the data sets used for regression analysis and fitting.

#### *Force-field based scoring functions*

Force-field based scoring functions calculate the sum of the ligand-protein interaction energy and the internal energy of the ligand. The affinity is estimated by force-field modeling of non-bonded interactions; summing the contributions from electrostatic interactions and van der Waals forces between the atoms of the ligand and the protein. Electrostatic terms are represented by the Coulombic potential-energy function (Equation 4) and van der Waals forces are calculated from the Lennard-Jones potential (Equation 5). As a last step, since the binding normally takes place in aqueous solution and water is not explicitly present in the model, the desolvation energies of the protein and the ligand are implemented by some scoring function. However force-field based scoring functions have major limitations due to complications in implementation; cut-off distances for non-bonded interactions have to be introduced, long-range effects should be accurately handled and solvation terms should be added. The GOLD<sup>41</sup>, G-Score<sup>47</sup>, D-Score<sup>47</sup> and DOCK scoring functions<sup>46</sup> are force-field based. AutoDock's scoring function<sup>78</sup>, on the other hand, combines empirical and force-field based terms.

#### *Knowledge-based scoring functions*

Based on statistical observations of intermolecular contacts in large 3D databases like Protein Data Bank<sup>79</sup>, knowledge-based scoring functions try to capture information hidden in the structural data of the protein-ligand complexes rather than the binding affinity used by empirical scoring functions. They try to reproduce the experimental structures by statistical analysis and simple atomic interaction-pair potentials, trying to implicitly capture the binding effects that are difficult to model explicitly. Knowledge-based methods use the structural information stored in databases of protein-ligand complexes to derive atom pair interaction potentials. These methods assume that the frequency of observing individual contacts reflect their energetic contribution to binding. When certain types of atoms interact more often than would be expected by a random distribution, they are likely to be energetically more favorable, thus contributing more to the binding affinity. PMFScore<sup>80,81,82</sup> and DrugScore's scoring functions<sup>83</sup> employ knowledge-based methods for binding affinity calculation.

However, all these methods are imperfect and often rank molecules poorly due to their intrinsic biases.<sup>84</sup> Knowledge-based and empirical methods are strongly influenced by the training sets used for fitting the function parameters or by the quality of the experimental data. It is also known that, regardless of the scoring function, larger molecules tend to produce better scores than smaller molecules simply because of the abundance of hypothetical interactions in the binding sites.<sup>37,85</sup> Accurate scoring and ranking is still a challenging problem, and different scoring functions can behave very differently in predicting the binding affinities of the same ligands to the same targets.<sup>55,73,84</sup> However, scoring functions are still needed for quick estimation of binding free energies of thousands of ligands in virtual screening studies.

### Consensus scoring

Consensus scoring is a well-known strategy used to improve the inaccurate ranking obtained with single scoring functions and to increase the number of true actives discovered. It involves evaluating a ligand-protein complex with several different scoring functions, which are then combined to reach a consensus conclusion. The logic behind this is that the combination of independent measurements leads to a value that is closer to the actual value and that the deficiencies of each scoring function may be compensated by the other scoring functions.<sup>86</sup> It has been reported that using different scoring functions to score ligand poses and deriving a consensus score that combines the outputs of these different scoring functions can improve the overall result in virtual screening studies.<sup>86,87,88</sup>

Consensus scoring strategies can be divided into three groups: rank-by-vote, rank-by-number and rank-by-rank.<sup>86</sup> Rank-by-vote gives a vote for each compound either according to its presence in the top  $n\%$  of the database for each scoring function, or if the score of the molecule is within the top  $n\%$  of the full range of scores obtained for the whole database. Rank-by-number averages the scores that are given by different scoring functions, allowing the introduction of weights. This is only applicable if all the scores are on the same scale, which can be achieved by normalization. Rank-by-rank approach simply averages the ranks that are output by different scoring functions.

### 1.7.2. FEP And TI

Free energy perturbation (FEP) theory, introduced by R.W. Zwanzig in 1954, is a statistical mechanics method used for computing free energy differences from molecular dynamics or Metropolis Monte Carlo simulations.<sup>68</sup> According to FEP theory, the free energy difference for going from state **0** to state **1** is obtained from:

$$\Delta F = F_1 - F_0 = -k_B T \ln \left\langle \exp \left( - \frac{H_1(\bar{X}) - H_0(\bar{X})}{k_B T} \right) \right\rangle_0$$

where  $T$  is the temperature,  $k_B$  is Boltzmann's constant, and the angle brackets denote an average over a simulation run for state **0**. The conformational average is taken from either molecular dynamics or Monte Carlo simulations.  $H_1$  and  $H_0$  are the energies of the system, computed using the Hamiltonians  $H_1$  (corresponding to state 1) and  $H_0$  (corresponding to state 0), respectively, and the coordinates of the particles ( $\bar{X}$ ) generated with the Hamiltonian  $H_0$ .

A problem of this method is that it requires the states **0** and **1** to be sufficiently close for their probability distributions to overlap, i.e. for configurations sampled with  $H_0$  to be also probable under  $H_1$ . This can be resolved by generating an appropriate number of intermediate steps between states **0** and **1**. For ligand-protein binding, state **0** typically refers to a ligand A (or a ligand A bound to a protein) and state **1** refers to a ligand B whose affinity for the protein's binding site is to be compared to that of A, making  $\Delta F$  the calculated free energy of perturbation of one ligand into the other in solution or in the active site. The relative binding free energy of ligands A and B can then be computed by using a thermodynamic cycle.

Binding free energies cannot be computed directly unless a proper reaction coordinate is known (which strictly is never the case).

An alternative derivation known as thermodynamic integration makes use of a coupling parameter ( $\lambda$ ) to bring the system from the initial to the final states:<sup>67</sup>

**Equation 7**

$$\Delta F = F_1 - F_0 = \int_0^1 \left\langle \frac{\partial H_\lambda(\vec{X})}{\partial \lambda} \right\rangle_\lambda d\lambda$$

In TI, simulations are performed at different  $\lambda$  values between 0 (state **0**) and 1 (state **1**), the ensemble average of the derivative of the Hamiltonian with respect to  $\lambda$  denoted by the angle brackets is calculated using the corresponding  $H_\lambda$  and the integration of the ensemble-average values is then evaluated numerically. As with the FEP method, a sufficient number of intermediate  $\lambda$  states needs to be simulated to obtain a smooth integration.

Even though both FEP and TI are very rigorous methods that enable explicit representation of water and protein flexibility, they are only applicable to systems where the structure of the protein and the approximate binding mode of the ligand are known. Both methods are limited by the high computational demands of a thorough sampling of configuration space, the accuracy of the force fields and the required proximity between the initial and final states, e.g. two similar ligands to be compared. Since FEP and TI are computationally expensive and stringent methods that require considerable initial setup calculations, they are still far away from large-scale implementation in virtual screening studies.

### 1.7.3. MM-PBSA

In the mid-range between scoring functions and FEP and TI, stands the Molecular Mechanics/Poisson-Boltzmann Surface Area approach. It is faster and computationally less demanding than FEP and TI, and at the same time more strict than the scoring functions.

In the MM-PBSA method, a molecular dynamics simulation of the ligand-target complex is done in a periodic box with explicit solvent and counterions, generating an ensemble of configurations of the system that is kept for post-processing.<sup>71,72</sup> Solvent and counterions are then removed from these configurations. Then, the binding free energies are calculated with Equation 8:<sup>72</sup>

**Equation 8**

$$G = \langle E_{MM} \rangle + G_{PBSA} - TS_{MM}$$

where  $G$  stands for the free enthalpy (free energy at constant pressure)  $\langle E_{MM} \rangle$  reflects the mean molecular mechanical energy,  $G_{PBSA}$  is the free enthalpy of (de)solvation, approximated by solving the Poisson-Boltzmann (PB) equation for the electrostatic part and including a solvent-accessible surface area (SA) term for the hydrophobic part, and  $TS$  stands for the entropic contribution and is taken from a quasi-harmonic or normal mode analysis of the trajectory. These contributions are averaged over the configurations extracted from the

molecular dynamics trajectory. Contrary to the previous methods, MM-PBSA only calculates the end states.

It has been reported that binding free energies calculated with the MM-PBSA approach shown statistically significant correlation to experimentally measured binding constants in various structure-based design studies with diverse target proteins.<sup>89,90</sup>

#### 1.7.4. LIE

The LIE method<sup>69</sup> is a semiempirical approach that is faster than FEP/TI, typically requiring a few hours per binding estimate, yet more accurate than empirical scoring functions. The approximations behind the LIE method, namely electrostatic linear response together with a nonpolar binding contribution that depends linearly on ligand size (representing hydrophobic effect, translational/rotational entropy loss, etc.), leads to a simple linear relation between the binding free energy and the difference in ligand-surrounding average potential energies between the bound and free states, i.e. between the compound immersed in water and enveloped in the binding pocket. These average energies are then calculated from sufficiently long molecular dynamics (MD) or Monte Carlo runs. The standard LIE method has been applied in combination with docking in lead optimization studies<sup>91,92,93</sup> and it has been suggested to have promising potential as a method that can be used at large scale.<sup>94</sup> The relationship between the ligand intermolecular interaction energies and the free energy of binding is given by the equation:

$$\Delta G_{bind} = \alpha \Delta \langle U_{l-s}^{vdw} \rangle + \beta \Delta \langle U_{l-s}^{el} \rangle + \gamma$$

**Equation 9**

While the  $\beta \Delta \langle U_{l-s}^{el} \rangle$  term represents the polar *l-s* contribution to the binding free energy and is based on a linear response approximation,  $\alpha \Delta \langle U_{l-s}^{vdw} \rangle + \gamma$  represents the *l-s* nonpolar binding contributions. The  $\beta$  coefficient can be derived from the linear response approximation, which predicts a value of 0.5.<sup>95</sup> However, based on rigorous FEP calculations in different solvents<sup>95,96</sup>, in the standard parameterization of the LIE method  $\beta$  is determined by the ligand's chemical groups, with values between 0.33 and 0.5. The value obtained for  $\alpha$  by fitting to experimental binding free energies is 0.18 and takes all size dependent contributions to binding into account, such as the hydrophobic effect and relative translational and rotational entropies as well as van der Waals interactions. The constant offset  $\gamma$  has been shown to correlate with the hydrophobicity of the binding site pocket<sup>96</sup> and is thus generally protein specific. Nevertheless, and for this reason, it can be safely ignored when screening molecules against one same protein.  $\langle U_{l-s}^{vdw} \rangle$  and  $\langle U_{l-s}^{el} \rangle$  are calculated, respectively, from the Lennard-Jones and electrostatic interactions between the ligand and its surrounding (*l-s*). These interactions are evaluated as energy averages (denoted by the angle brackets) from separate MD simulations of the free and bound states of the ligand (solvated in water and bound to the protein in solution, respectively). The difference ( $\Delta$ ) between the averages for the two states is then calculated.<sup>69</sup>

## *Chapter 2*

### *Objectives*



## ***2. Objectives***

In light of the information provided in the previous introductory chapter, the objectives of this thesis can be listed as follows:

1. To design an integrative approach for the computational discovery of drug-like small molecules with affinity for a target protein, based on available methodologies. This objective can be divided into smaller objectives:

- to find a robust docking protocol that can be applied to any drug target.
- to integrate pharmacophore filtering to high-throughput docking and to investigate the possible advantages and disadvantages of using pharmacophore filters.
- to find out the most convenient methodology for consensus scoring.
- to apply large-scale molecular dynamics simulations for late-stage improvement of ligand ranking.

2. to devise an evolving methodology that enables addition of new approaches at different stages; starting as a simple methodology based on high-throughput docking with consensus scoring, then getting to a more complicated hybrid methodology that integrates pharmacophore filtering and LIE simulations.

3. to apply the approach developed in the previous objectives to find candidate ligands for three target proteins: human T-protein, human bleomycin hydrolase and human acid  $\beta$ -glucosidase.





## *Chapter 3*

### *The Three-Step Docking Procedure Versus The Hybrid Procedure*



## ***3. The Three-Step Docking Procedure Versus The Hybrid Procedure***

### ***3.1. Summary***

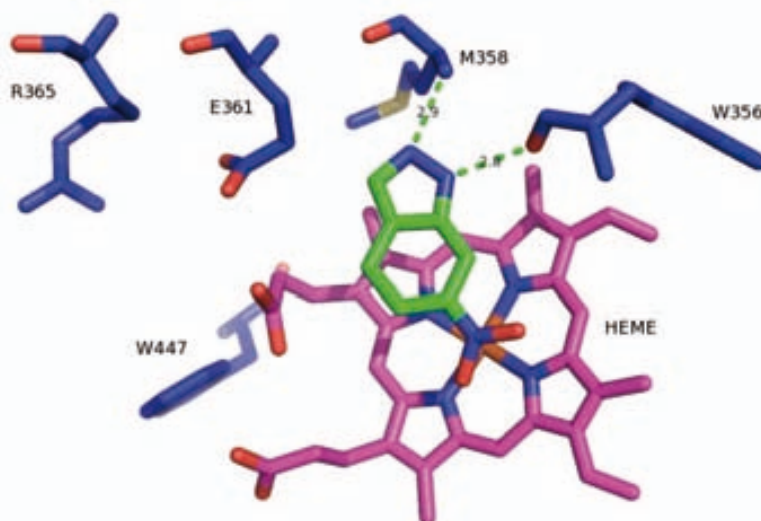
For a very complex problem like drug discovery, using a single straightforward computational method or tool is not likely to give reasonable and reliable results. As mentioned earlier, the accuracy of docking and scoring is highly dependent on the target; a docking tool producing ligand binding modes very close to the experimental mode for a specific target may not be as successful for another target.<sup>53,54</sup> With the aim of developing a molecular docking based automated procedure that can be applied to any target and that produces relevant binding modes, a three-step docking procedure was developed. This procedure consists of a very trivial size filter to exclude molecules too large for the binding site, a three-step docking with AutoDock 4.0 and ranking with consensus scoring at the end of each docking step. The docking parameters were calibrated by docking a sample set from a small molecule library to a test protein, human endothelial nitric-oxide synthase (eNOS) enzyme, and the parameters were validated with a benchmark protein with known active ligands, human checkpoint kinase 1. The success of the three-step docking procedure was also compared with a hybrid procedure, which is a combination of pharmacophore filtering, automated docking with AutoDock Vina and consensus scoring. Both procedures were applied to human checkpoint kinase 1 to find out the optimal strategy.

### ***3.2. Biological Background***

#### ***3.2.1. Sample Protein: Human Endothelial Nitric-Oxide Synthase (eNOS)***

Nitric oxide synthases (NOSs) are a family of enzymes that catalyze the production of nitric oxide (NO) from L-arginine.<sup>97</sup> Nitric-oxide is an important cellular signaling molecule that has a vital role in different biological mechanisms such as controlling vascular tone (hence blood pressure), airway tone, insulin secretion, and peristalsis.<sup>98</sup> It is also involved in the development of the nervous system and in angiogenesis.<sup>98</sup> Nitric oxide signaling is mediated in mammals by three isoenzymes, eNOS (endothelial NOS), nNOS (neuronal NOS) and iNOS (inducible NOS) involved in immune response.<sup>97,99,100</sup> All NOSs are homodimeric enzymes consisting of two conserved modules:<sup>88</sup> an electron-supplying reductase module and a catalytic oxygenase module (NOS<sub>ox</sub>).<sup>101</sup> Endothelial NOS (eNOS) is a nitric oxide synthase that generates NO in blood vessels and is involved in regulating the vascular tone by inhibiting smooth muscle contraction and platelet aggregation.<sup>102</sup> Variations in the gene coding eNOS are associated with susceptibility to coronary spasm, and a polymorphism in the gene has been shown to be connected to Alzheimer's disease.<sup>103,104</sup>

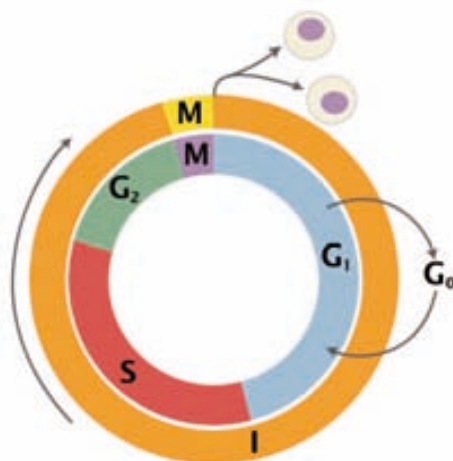
To calibrate the docking parameters for the three-step docking procedure of a library of nearly 2,000,000 ligands we selected eNOS, which has known ligands and active site. The three-dimensional structure for the receptor was obtained from the Protein Data Bank (PDB) and corresponds to entry 1M9M.<sup>105</sup> This structure has a heme group as one of the cofactors along with a zinc ion and a known ligand (6-nitroindazole, 6NI) bound (Figure 4).



**Figure 4:** Structure of eNOS<sub>ox</sub> binding site co-crystallized with 6-nitroindazole (active site inhibitor orientation). Hydrogen bonds with the Met358 amide nitrogen and the Trp356 carbonyl oxygen are shown as green dashed lines.

### 3.2.2. Benchmark Protein: Human Checkpoint Kinase 1

The cell cycle is the series of events that takes place in a cell leading to its division and duplication and consists of four distinct phases: G<sub>1</sub> phase, S phase (synthesis), G<sub>2</sub> phase and M phase (mitosis).<sup>106</sup> G<sub>1</sub>, S and G<sub>2</sub> phases together are called interphase, the phase of the cell cycle in which the cell spends the majority of its time and prepares for cell division. Cells that have temporarily or reversibly stopped dividing enter a state of resting called G<sub>0</sub> phase (Figure 5). Activation of each step depends on the proper progression and completion of the previous phase, with two checkpoints throughout the cell cycle to ensure completion of the preparation for cell division.<sup>107</sup> The first checkpoint is located at the end of the G<sub>1</sub> phase, before entry into the S phase (G<sub>1</sub>/S checkpoint), making the key decision of whether the cell should divide, delay division or enter a resting stage. The second checkpoint is located at the end of the G<sub>2</sub> phase (G<sub>2</sub>/M checkpoint), deciding the initiation of mitosis.

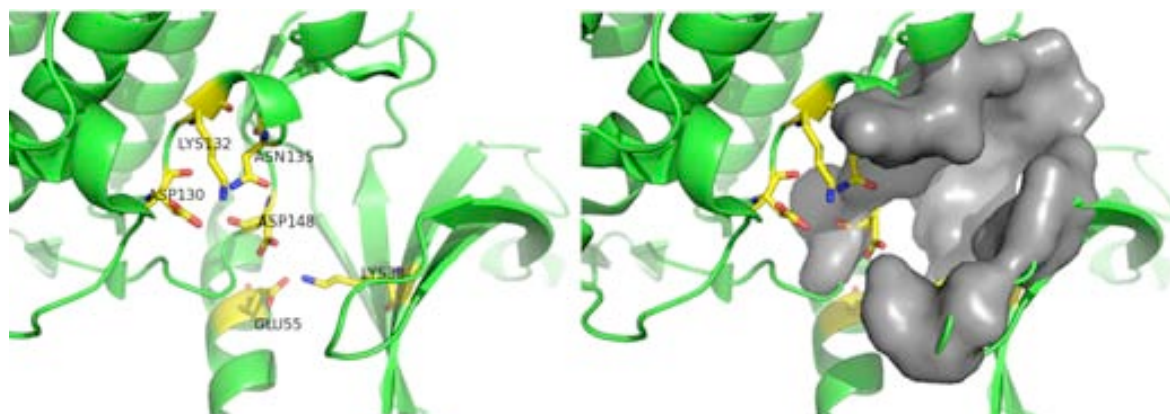


**Figure 5:** Schematic representation of the cell cycle. Outer ring: I: Interphase, M: Mitosis; inner ring: M: Mitosis, G<sub>1</sub>: Gap 1, G<sub>2</sub>: Gap 2, S: Synthesis; not in ring: G<sub>0</sub> = Gap 0/Resting. (Figure taken from [http://en.wikipedia.org/wiki/Cell\\_cycle](http://en.wikipedia.org/wiki/Cell_cycle))

The understanding of the cell cycle provides the identification of potential targets that may become key elements for the development of therapeutics against cancer.<sup>108</sup> An important stage in the cell cycle is the G<sub>2</sub>/M checkpoint, which ensures that cells don't initiate mitosis before they have a chance to repair damaged DNA after replication.<sup>109</sup> If a cell with a damaged DNA passes through a defective G<sub>2</sub>/M checkpoint without DNA repair and enters mitosis, it ultimately leads to cell death.<sup>110</sup> In many tumor cells, the first checkpoint at G<sub>1</sub>/S is impaired and tumor cells are not arrested at this stage, leaving the G<sub>2</sub>/M checkpoint as the only control. If the G<sub>2</sub>/M checkpoint is abrogated in these tumor cells, they will enter mitosis prematurely with DNA damage resulting in cell death.<sup>111</sup> Therefore, inhibition of the G<sub>2</sub>/M checkpoint leads to selective sensitivity of cancer cells and, hence, small molecule inhibitors of this checkpoint may have potential use as sensitizing agents for cancer therapy.<sup>112</sup>

To test the three-step docking procedure and to compare it with a hybrid method, we used the human checkpoint kinase 1<sup>113,114</sup> (Chk1), which arrests cells with DNA damage at the G<sub>2</sub>/M checkpoint and prevents cell division. Chk1 has emerged as a promising target for the design of small molecule inhibitors and has been studied extensively in therapeutic research for cancer.<sup>111,112,115</sup> The motivation behind choosing Chk1 as the validation protein is that it has many known inhibitors available with experimentally proven activity values.

The structure with PDB id 1IA8<sup>116</sup>, which corresponds to apo human Chk1, was used. The actual substrate of Chk1 is ATP and therefore the ATP binding site of 1IA8 has been determined to be the active site for finding competitive inhibitors. The catalytic site residues (Figure 6) are Lys38, Glu55, Asp130, Lys132, Asn135 and Asp148.<sup>116</sup> However, the residues involved in ligand binding can be extended to include a buried pocket containing Leu15, Gly16, Gly18, Tyr20, Val23, Ala36, Val68, Leu82, Leu84, Glu85, Tyr86, Cys87, Ser88, Gly90, Glu91, Glu134, Leu137, Ser147 and Phe149.<sup>111,112,115,117</sup>



**Figure 6:** Catalytic residues and binding site of Chk1. Left side of the figure shows the catalytic site residues of Chk1. On the right side, residues that are part of the binding site are displayed as a grey surface to show the pocket-like structure that creates the binding site.

### ***3.3. Common Methods For Three-Step Docking And The Hybrid Procedure***

The three-step docking approach and the hybrid approach basically use the same outline, however with different applications. They both consist of a primary filter to reduce the number of molecules, followed by docking and consensus scoring. The same small molecule and protein files are used as input for both methods. However, they differ at the application of the filter and the docking steps. While the three-step docking approach uses a size-filter and a three-step docking procedure with AutoDock 4.0, the hybrid method works with a pharmacophore filter and a single step docking with AutoDock Vina. This section focuses on the common parts of these two approaches; the preparation of the small molecule libraries and the proteins, consensus scoring and the selection of known inhibitors of Chk1. Section 3.4 presents the methodologies used only by the three-step docking approach while Section 3.5 elaborates on the details specific for the hybrid approach.

#### ***3.3.1. Preparation Of The Small Molecule Libraries And The Proteins***

##### *The preparation of the small molecule library*

The small-molecule library (VSL-1) used in this part of the work is based on the compilation of compounds found in the Chemical DataBase Manager (CDBM), built by the group of Dr. Xavier Barril at the Department of Physical Chemistry of Universitat de Barcelona and containing commercially available compounds from several vendors. In CDBM, molecular configurations (states) and three-dimensional conformations are generated with LigPrep<sup>118</sup>, enumerating tautomers, ionization states and, when the chirality is not specified, enantiomers. The generated configurations are then minimized using the OPLS force-field<sup>119</sup>. To this raw library, filters such as the Lipinski Rules<sup>120</sup>, Veber Rules<sup>121</sup> and not having reactive moieties or more than 4 states (tautomeric, ionization and enantiomeric) are applied. This filtered virtual screening library contains 1,961,165 entries stored in 393 multi-SD files each containing 5000 molecules. From these SD files, we generated UNITY<sup>34</sup> databases for pharmacophore search for the hybrid procedure, PDBQT<sup>122</sup> files (with MGLTools 1.5.4<sup>123</sup>)

for docking with AutoDock 4.0 for the three-step-docking procedure and AutoDock Vina for the hybrid procedure, and MOL2 files<sup>124</sup> for consensus scoring with CSCORE<sup>87</sup>.

VSL-1 was converted into several different formats required by the different types of software used in this work. Using the dbimport function of UNITY<sup>34</sup> with the parameters +perceive\_chiral\_c, +perceive\_chiral\_np, +perceive\_stereo\_bond, and -prescan 0, the UNITY database was created from the SD files. AutoDock 4.0 and AutoDock Vina require the compounds in PDBQT format<sup>122</sup>. The conversion of the whole library from SD format to individual PDBQT files for each compound was done in two steps. First, the multi-structure SD files were converted into multi-structure MOL2 files<sup>124</sup> using UNITY's dbtranslate function. The resulting multi-structure MOL2 files were split into individual files for each compound, which were then converted to PDBQT format using AutoDock MGLTools 1.5.4<sup>123</sup>. However, during the conversion to PDBQT sometimes problematic "empty branches", i.e., a BRANCH statement immediately followed by an ENDBRANCH statement in the PDBQT file, were created. All PDBQT files were scanned for such patterns, and were fixed with a program that was developed only for this task. This program checked the atom order and coordinates from the PDBQT and corresponding MOL2 files, generating PDBQT files with correctly defined atom and bond types.

#### *The preparation of the protein files*

The PDB entry 1M9M has a heme group as one of the cofactors along with a zinc ion and a known ligand (6-nitroindazole, 6NI) bound. All solvent molecules and 6NI were removed from the coordinate file, leaving only the protein and the heme group. The heme group wasn't removed because it is a part of the binding site and the known ligand; 6NI sits on it in the bound state. The docking experiments were done only on chain A of 1M9M with the heme group bound. For AutoDock's file format PDBQT, Gasteiger charges<sup>125</sup> were added and non-polar hydrogens were merged to united carbon atoms. For the three-step-docking procedure the structure of human eNOS was only used to determine the parameters to be used for high-throughput docking for Chk1.

For the benchmark, the crystal structure of apo human Chk1 with PDB id 1IA8 was used. For the three-step-docking procedure with AutoDock 4.0 and the AutoDock Vina docking part of the hybrid procedure, the PDBQT file for Chk1 was prepared in the same way as eNOS: water molecules were removed, non-polar hydrogens were included in united atoms and Gasteiger charges were added. However, for the consensus scoring and pharmacophore searches, an alternative treatment was needed. This was done with the Biopolymer Structure Preparation tool of SYBYL-X<sup>34</sup>: water molecules were removed, hydrogens and charges (AMBER7\_F99 charge set) were added and a MOL2 file was created as input for pharmacophore searching and consensus scoring with CSCORE.

### **3.3.2. Collection Of The Known Actives**

To test the efficiency of the three-step docking and the hybrid methods and to compare them, we collected 21 known actives of Chk1 and added them to the VSL-1 library to see if these two methods would be able to pick the known ligands among the other molecules (Table 2 and Figure 7). Three of the known ligands are natural products such as staurosporine<sup>126</sup> and its derivatives and the rest were ligands that have been identified with virtual screening methods.

Molecules **A1** to **A6** and **B1** to **B3** were already present in the library, but the rest of the molecule structures were downloaded from PubChem<sup>127</sup> and treated in the same way as the other compounds in VSL-1, as explained in section 3.3.1. Molecule **E1** is the natural product staurosporine, and **E2** and **E3** are derivatives of it.

Molecule Name	Reference	In VSL-1
A1, A2, A3, A4, A5, A6	Foloppe et. al. <sup>115</sup>	Yes
B1, B2, B3	Foloppe et. al. <sup>111</sup>	Yes
B4, B5, B6, B7	Foloppe et. al. <sup>111</sup>	No
C1, C2, C3	Foloppe et. al. <sup>117</sup>	No
D1, D2	Lyne et. al. <sup>112</sup>	No
E1, E2, E3	Zhao et. al. <sup>126</sup>	No

**Table 2:** The references for the molecules chosen for testing the methods.



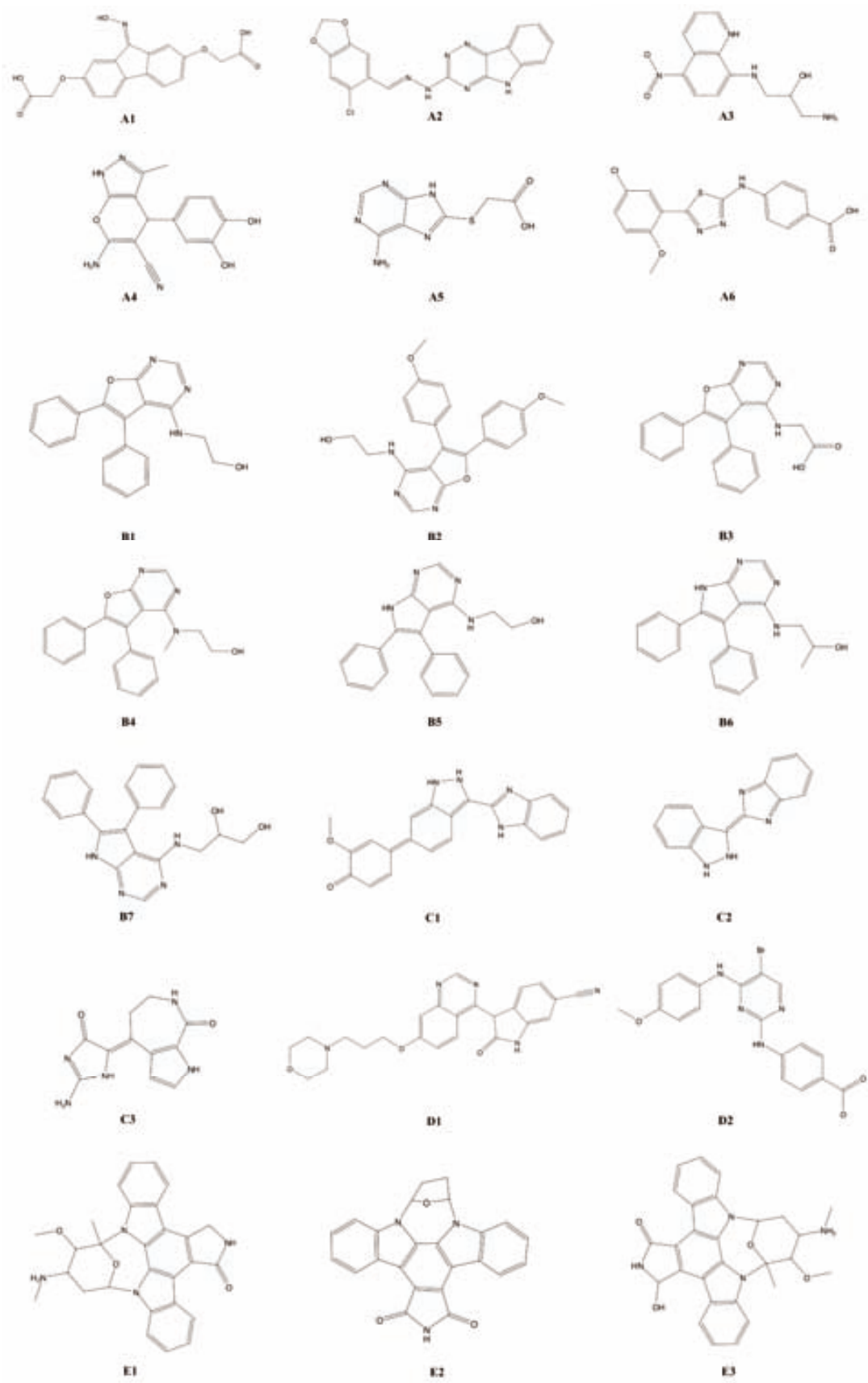


Figure 7: Known actives of Chk1.

### 3.3.3. Consensus Scoring

In this study, a given ligand-receptor complex produced by AutoDock 4.0 or AutoDock Vina was re-scored using D-Score, PMFScore, ChemScore and G-Score with the CSCORE program of SYBYL-X. CSCORE program requires MOL2 files, however AutoDock 4.0 and AutoDock Vina output PDBQT files. Direct file format conversion is problematic because of different atom type identifiers and atom orderings used by these two file formats. Therefore, we wrote a small script that orders the atoms of the ligands as in the MOL2 files of the ligand library and takes the corresponding docked atomic coordinates from the PDBQT files created by AutoDock. This way, we created proper MOL2 files with correct atom types, bonds and bond types, avoiding error-prone format conversion. Hydrogen atoms that were missing in the created MOL2 files because of AutoDock's united-atom approach were added using SYBYL's FILLVALENCE function. The four scores produced by CSCORE were combined with the score obtained with docking. However, instead of CSCORE's rank-by-vote consensus scoring approach, a modified rank-by-number approach was used to combine these five scoring functions. For each scoring function, all scores were normalized to values between 0 and 1 (0 representing the most favorable compound, 1 representing the least favorable one by each scoring function), such that they were on the same scale and therefore comparable. Two different normalization procedures were implemented. The first one was applied to all scores of all docked ligands. The second one was applied after truncating the 0.5% most poorly scoring molecules for each scoring function by directly assigning 1 as their normalized scores. In other words, in the second procedure 99.5% of the scores were normalized after the deletion of the poorly scoring 0.5% part.

The normalized score of compound  $i$  with scoring function  $F$ ,  $S_{cut-off,F}(i)$ , was thus defined as,

Equation 10

$$S_{cut-off,F}(i) = \min\left(1, \frac{E_F(i) - E_{\min,F}}{E_{cut-off,F} - E_{\min,F}}\right)$$

where  $E_F(i)$  is the score of compound  $i$  (the best pose given by docking) with scoring function  $F$ ,  $E_{\min,F}$  is the lowest (best) score obtained with the scoring function  $F$  and  $E_{cut-off,F}$  is the value of the first score falling above the  $cut-off\%$  of the scores obtained with function  $F$ . This equation defines both normalization procedures, with or without truncation. In the first normalization procedure, applied to all scores of all docked ligands, setting the truncation cut-off to 100% implies no truncation. As already mentioned, in the second procedure the cut-off was set to 99.5%.

To obtain a consensus score over all scoring functions we summed all normalized scores for a given compound, yielding the "normalized consensus score". The normalized consensus score of compound  $i$ ,  $NCS_{cut-off}(i)$ , is thus defined as Equation 11:

Equation 11

$$NCS_{cut-off}(i) = \sum S_{cut-off,F}(i)$$

The compound with the smallest normalized consensus score is taken as the best binder.

### ***3.4. Methods For The Three-Step Docking Approach***

This section elaborates on the methodologies used only in the three-step docking approach. The docking part of the approach consists of three runs; i) a first rapid docking step, ii) a second docking step, the exhaustiveness of which is decided according to the number of rotatable bonds of the docked molecule, and iii) a very thorough third step. At the end of each step, the molecules are ranked according to their consensus scores and only a fraction goes through the next step of dockings.

#### ***3.4.1. Parameter Calibration For Docking With The Sample Protein***

The process of docking with AutoDock 4.0 can be divided roughly into three main steps: preparation of grid files and calculation of atom types for each ligand, the assignment of genetic algorithm docking parameters, and the docking simulation itself. The results depend on an appropriate choice of parameters for the docking procedure. With the aim of automating the parameter assignment for the second step of this approach according to the number of rotatable bonds of the ligand to be docked, we decided to choose human eNOS as sample protein.

##### *Preparation of grid files and calculation of atom types for each ligand*

In the AutoDock 4.0 docking algorithm, the interaction between a putative (probe) ligand atom and the receptor's binding site is pre-calculated and given as a set of grid-based potential energy files called 'gridmaps'. Each of these gridmap files is specific for a given atom type. If an atom type in the ligand is absent in the gridmap file, it is not included in the interaction calculation with the receptor atoms. Therefore, the gridmap files should be prepared specific for each ligand, covering all atom types. However, when screening a library that contains almost two million molecules, calculating the gridmap files for each ligand is computationally very costly. In addition, almost all ligands contain the atom types C, N and O, and calculating gridmaps for these atoms for each ligand is therefore redundant. To overcome this inefficiency and redundancy, and to decrease the time needed, we calculated the interaction grids for all possible atom types and used the potential gridmap files created uniquely for the receptor. To this end, a grid with 58 x 54 x 60 points and a default spacing of 0.375 Å was centered at the binding site of the chain A of the protein.

##### *Assignment of genetic algorithm docking parameters using the sample protein*

In the Lamarckian genetic algorithm implemented in AutoDock, the key parameters are the number of docking runs, the maximum number of generations and the maximum number of energy evaluations. As these numbers increase, the accuracy of the docking procedure increases, as does the computational effort.

To find out the optimal parameters for each ligand, two docking experiments were done with eNOS. First, all molecules in the library were docked to the protein and the results were ranked according to their binding free energy calculated by AutoDock 4.0. In a second step, 1000 molecules were chosen as a sample set and re-docked to the receptor protein with different combinations of docking parameters to find out which parameter sets to be used according to the number of rotatable bonds of the small molecules.

In the first docking experiment, all the molecules in VSL-1 were docked within the grid centered on the binding site of eNOS with 10 docking runs per compound and the maximum number of generations and energy evaluations set to 27,000 and 250,000, respectively. These parameters correspond to a very rapid virtual screening done to differentiate possible hits with negative binding free energies from molecules that are not viable. At the end of the docking simulations, the minimum-energy pose of each docked ligand was taken and ranked according to its calculated binding free energy.

From the ranked ligands, we chose a sample set of 1000 molecules. To define this set, we first divided the energy space spanned by all molecules into two: molecules with positive binding free energy and molecules with negative binding free energy. For the positive part, the first 100 and the last 100 molecules were taken. The negative energy space was divided into eight parts and 100 molecules from each part were selected. Thus, the test set contained 1000 molecules: 200 molecules from the positive energy space and 800 molecules from the negative energy space. This set was chosen with the aim to span the energy space obtained with the whole library, thus being a representative set for the whole library. The second docking experiment was done only with this sample set.

In the second round of dockings for parameter assignment, sets of different values for each parameter of the genetic algorithm were used (Table 3). The values were 10 and 100 for the number of docking runs, 250,000, 500,000, 1,000,000 and 2,500,000 for the maximum number of energy evaluations, and 10,000 and 27,000 for the maximum number of generations. Each of the 1000 selected molecules was docked with all possible combinations of these parameters, resulting in 16 different docking simulations per ligand.

ga_run	num_eval	num_gen
10	250.000	10.000
100	500.000	27.000
	1.000.000	
	2.500.000	

**Table 3:** Parameter values used for docking of the test set of molecules to human eNOS.

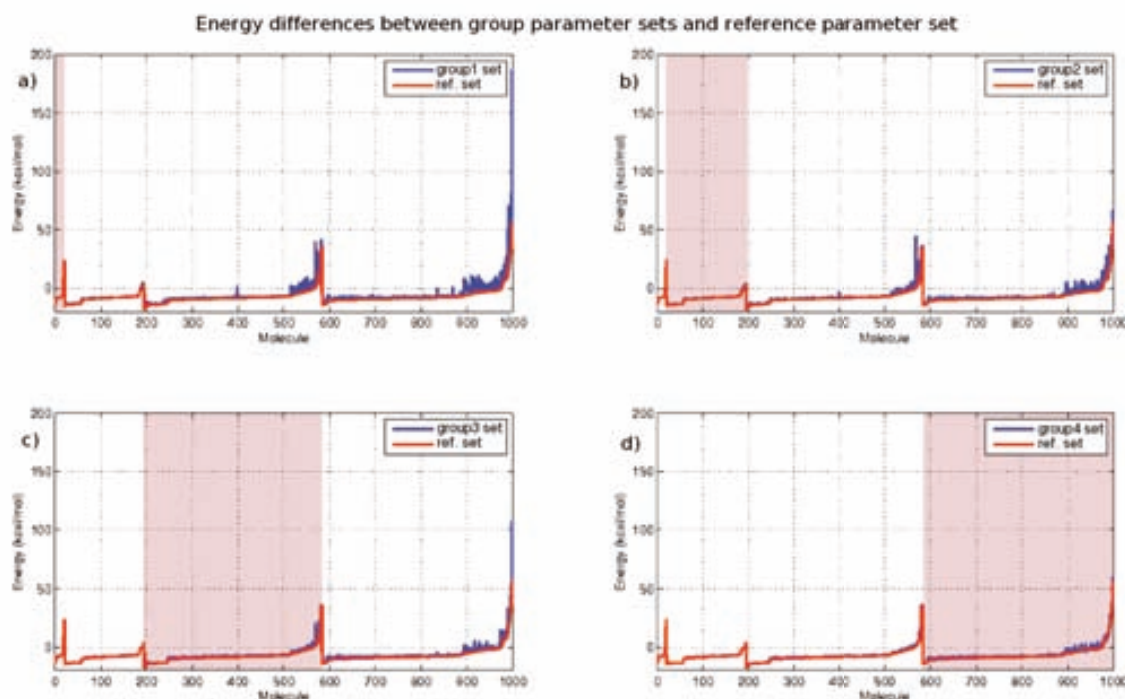
We divided the test set into four groups according to the number of rotatable bonds of the molecules (Table 4) and calculated the error in the resulting binding free energy for each group and for the different parameter sets, taking as reference value the one obtained with the most thorough parameter combination. The aim of this step is to determine the correlation between the number of rotatable bonds of a molecule and the parameters to be used for an efficient docking that molecule. Docking each one of the molecules in the test set to the target protein's binding site using 16 different parameter sets revealed, as expected, that molecules with fewer number of rotatable bonds can be correctly docked with a faster docking, i.e., using smaller values for the maximum number of energy evaluations and the number of docking runs. On the other end, molecules with several rotatable bonds require relatively large values for the maximum number of energy calculations and the number of docking runs.

Molecule Groups	Description
Group 1	Molecules with 0 or 1 rotatable bond
Group 2	Molecules with 2 or 3 rotatable bonds
Group 3	Molecules with 4 or 5 rotatable bonds
Group 4	Molecules with more than 5 rotatable bonds

**Table 4:** Molecule groups divided according to the number of rotatable bonds.

We assumed the results with the most thorough parameter combination (ga\_run:100, num\_evals:2.500.000 and num\_gen:27.000) were the most accurate and calculated their differences in energy with the remaining 15 docking results using alternative parameter sets, to obtain an error value for each parameter set and molecule group. According to this, the final parameter set-molecule group coupling was (Figure 8):

- For molecules with 0 or 1 rotatable bond (Group 1): (ga\_run=10, num\_eval=250.000, num\_gen=27.000).
- For molecules with 2 or 3 rotatable bonds: (ga\_run=10, num\_eval=500.000, num\_gen=27.000).
- For molecules with 4 or 5 rotatable bonds: (ga\_run=10, num\_eval=1.000.000, num\_gen=27.000).
- For molecules with more than 5 rotatable bonds: (ga\_run=10, num\_eval=2.500.000, num\_gen=27.000).



**Figure 8:** The energy differences between group parameter sets and the reference set. **a)** Energy difference between docking simulations with the Group1 parameter set (ga\_run: 10, num\_eval: 250.000, num\_gen: 27.000) and the most thorough parameter set (ga\_run: 100, num\_eval:2.500.000, num\_gen: 27.000). The docking done with the Group1 parameter set is approximately 100 times faster than the docking done with the thorough parameter set. The highlighted part contains molecules from Group1. In this figure, it is seen that using this parameter set would cause the molecules from groups 3 and 4 to be poorly scored. **b)** Energy difference between docking simulations with the Group2 parameter set (ga\_run: 10, num\_eval: 500.000, num\_gen: 27.000) and the most thorough parameter set. The docking done with the Group2 parameter set is approximately 50 times faster than the docking done with the thorough parameter set. The highlighted part contains molecules from Group2. **c)** Energy difference between docking simulations with the Group3 parameter set (ga\_run: 10, num\_eval: 1.000.000, num\_gen: 27.000) and the most thorough parameter set. The docking done with the Group3 parameter set is approximately 25 times faster than the docking done with the thorough parameter set. The highlighted part contains molecules from Group3. **d)** Energy difference between docking simulations with the Group4 parameter set (ga\_run: 10, num\_eval: 2.500.000, num\_gen: 27.000) and the most thorough parameter set. The docking done with the Group4 parameter set is approximately 10 times faster than the docking done with the thorough parameter set. The highlighted part contains molecules from Group4.

### ***3.4.2. The Size Filter***

Before the computationally expensive docking experiments, it is a good idea to filter out in a fast and simple way the molecules that cannot fit in the binding site. Thus, we filtered our library of small molecules according to size; we calculated the diameter of the binding site and filtered out the molecules whose length is larger than 1.5 times the binding site diameter. The cut-off size corresponding to the benchmark protein (Chk1) was calculated to be 19 Å and applied to the ligand configurations present in the library (excluding non-polar hydrogens). However, since the binding site of Chk1 is rather large (Figure 6), the size filter couldn't reduce the number of molecules effectively in this case. The VSL-1 library was reduced to 1,918,830 molecules from 1,961,165 after the application of this size filter.

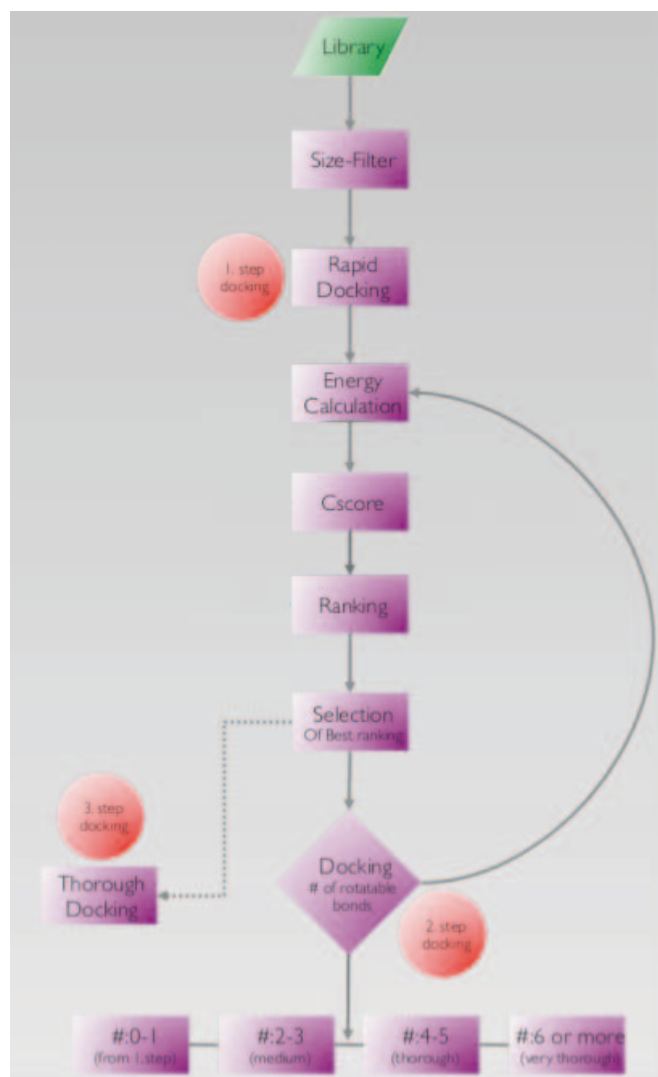
### ***3.4.3. Docking Experiments For The Three-Step Docking Approach***

All three docking steps took place on a grid centered on the ligand binding site of Chk1 with 58 x 60 x 48 points and 0,375 Å spacing, prepared with MGLTools 1.5.4 of AutoDock 4.0. The grid was made large enough to include the binding site residues and the buried pocket. The same gridmaps were used for all docking steps done with AutoDock 4.0.

We integrated parameter-set and molecule-group coupling into an automated docking process, namely the three-step docking approach (Figure 9). Prior to the computation-intensive docking experiments, we filtered the library of small molecules according to size, by excluding the molecules whose lengths were more than 1.5 times the binding site diameter. Later, this filtered library was docked to the receptor protein with the same parameter set (ga\_run:10, num\_eval:250.000 and num\_gen: 27.000) for all groups of molecules. The parameter set for the first docking step corresponds to a very rapid docking, performed to differentiate the molecules that score really poorly and exclude them from further steps.

All the molecules from the first step were ranked according to their normalized consensus scores and the best scoring 60000 were selected for the second step of dockings. This was performed with parameters determined according to the number of rotatable bonds of the compounds, calibrated previously with the parameter set testing (see section 3.4.1). The only differing parameter is num\_evals; ga\_run and num\_gen are fixed to 10 and 27.000, respectively.

For the third step, again the conformations resulting from the second step were scored with five different scoring functions and their normalized consensus scores were calculated. The best 1000 were chosen for a very thorough docking procedure with parameters ga\_run:100, num\_evals:2.500.000 and num\_gen: 27.000.



**Figure 9:** The flowchart summarizing the three-step docking algorithm. First the library of small molecules is filtered according to size and then all the molecules that pass the filter are docked to the protein with the rapid first docking step. Then energies of the poses are calculated by different scoring functions and the different scores are combined to a single score with consensus scoring. The molecules ranking at the top after consensus scoring are chosen for a second step of docking, however this time the exhaustiveness of docking is decided based on the number of rotatable bonds of the molecule. At the end of the second round of docking, energy calculations and consensus scoring are done again to select the final set of molecules that will undergo a third round of docking, which is done very thoroughly. A final step of energy calculation and consensus scoring is done to find the final ranking of the molecules.

### 3.5. Methods For The Hybrid Approach

The hybrid approach integrates pharmacophore filtering instead of a size filter and reduces the number of docking steps to only one because pharmacophore filtering prunes the libraries to a larger extent than the size filter (Figure 10). For the hybrid approach, the same molecule libraries and the structure of 1IA8 used for the three-step docking method were used (section 3.3.1). However, for the docking part of the hybrid approach AutoDock Vina was used instead of AutoDock 4.0. Consensus scoring was done in the same manner as for the three-

step docking with the only difference being that the binding free energies of the ligands are calculated with AutoDock Vina instead of AutoDock 4.0.



**Figure 10:** The flowchart summarizing the hybrid method. First the library of small molecules is filtered with a pharmacophore filter and then all the molecules that pass the filter are docked to the protein with AutoDock Vina. Then energies of the poses are calculated by different scoring functions and the different scores are combined to a single score with consensus scoring. The molecules are ranked according to this normalized consensus score, and the top-ranking molecules are chosen to be the leads.

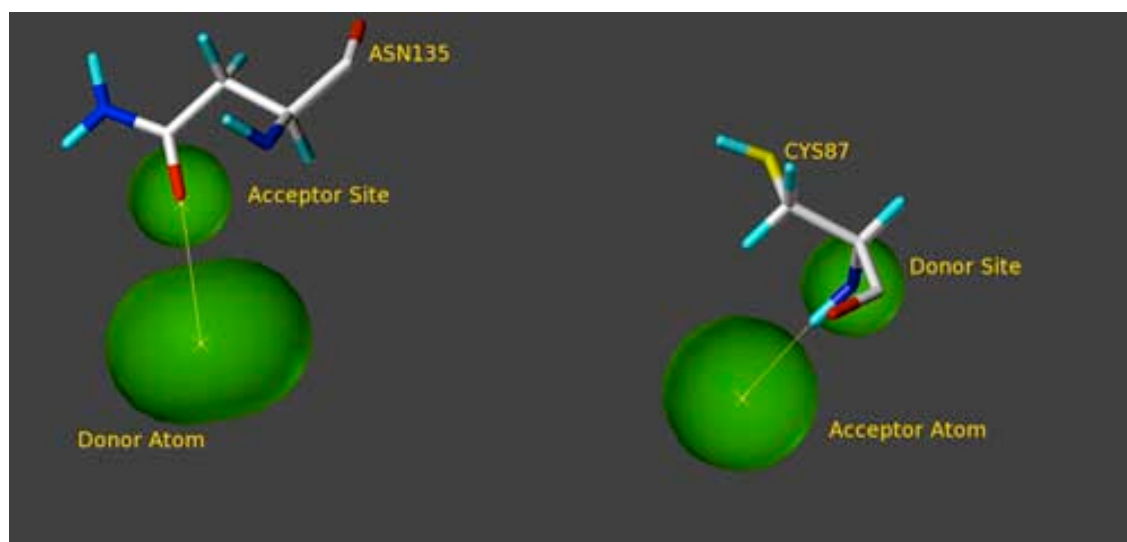
### 3.5.1. Pharmacophore Filtering

One of the differing parts between the three-step docking and the hybrid approaches is the use of pharmacophore filters and searches. We used two separate pharmacophore filters to reduce the number of molecules that would go under docking. Pharmacophore filters are more specific than the size filters and can weed out molecules that are not chemically or sterically compatible with the binding site.

The first pharmacophore filter was prepared with features only from the binding site of Chk1 (Figure 11). A hydrogen bond donor site with a van der Waals tolerance of 1 Å was placed on the nitrogen of the Cys87, requiring a hydrogen bond acceptor atom with a 1.5 Å van der

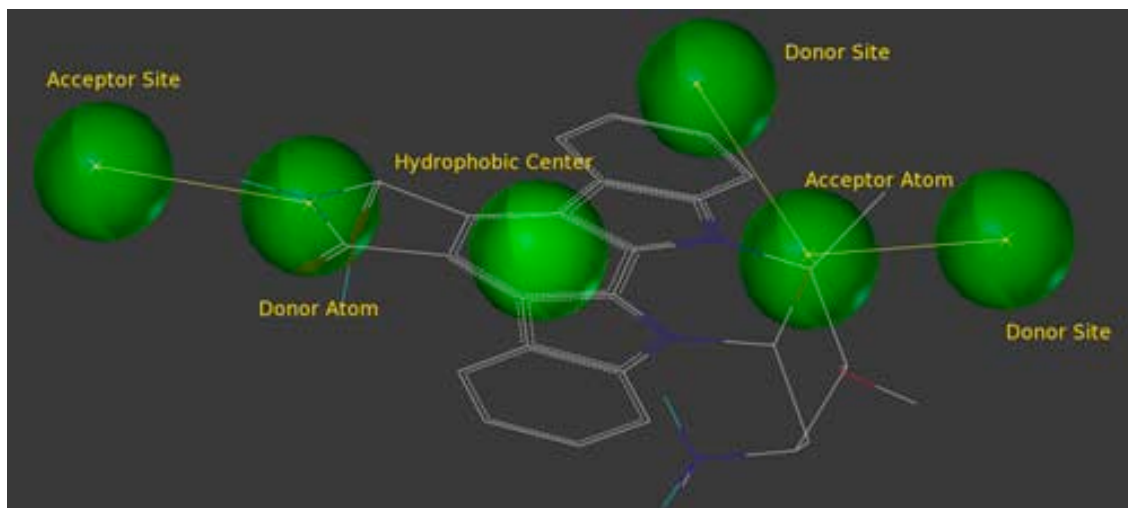


Waals tolerance on the candidate ligand. The other pharmacophore feature added was a hydrogen bond acceptor site with a 1 Å tolerance on the side chain oxygen of Asn135. This acceptor site required a hydrogen bond donor atom with a 1.5 Å tolerance on the candidate ligand. The remaining binding site residues were added as excluded volumes with a scaling factor of 0.25 for van der Waals atom radii. 4 of the 21 known ligands passed this filter: **A4**, **B4**, **D1** and **D2**. The library was reduced to 65671 molecules from 1.96 million with this filter.



**Figure 11:** The pharmacophore filter created from Chk1 binding site residues Cys87 and Asn135 (excluded volumes not shown).

The second pharmacophore filter was prepared with features taken from 3 of the known ligands—the natural products **E1**, **E2** and **E3**—and the receptor binding site constraint (Figure 12). The docked conformations of these three ligands output by AutoDock Vina were superimposed and the properties common to all three were chosen to be present in the pharmacophore filter. A hydrophobic and aromatic center was defined since all three molecules have 5 rings in a planar structure. The van der Waals tolerance for this hydrophobic and aromatic center was 0.5 Å. A hydrogen bond donor atom feature with 0.5 Å tolerance, needing a hydrogen bond acceptor site with 0.5 Å tolerance on the protein, was also located in the query. Finally, a hydrogen bond acceptor atom with 0.5 Å tolerance, binding to the protein hydrogen bond donor site with 0.5 Å tolerance, was also defined. The only pharmacophoric features added on the sole basis of the binding site characteristics were the excluded volumes, derived from the binding site residues and scaled with a factor 0.25. The library was reduced to 172339 molecules while 5 of 21 known ligands passed the filter: **A3**, **A5**, **E1**, **E2** and **E3**.



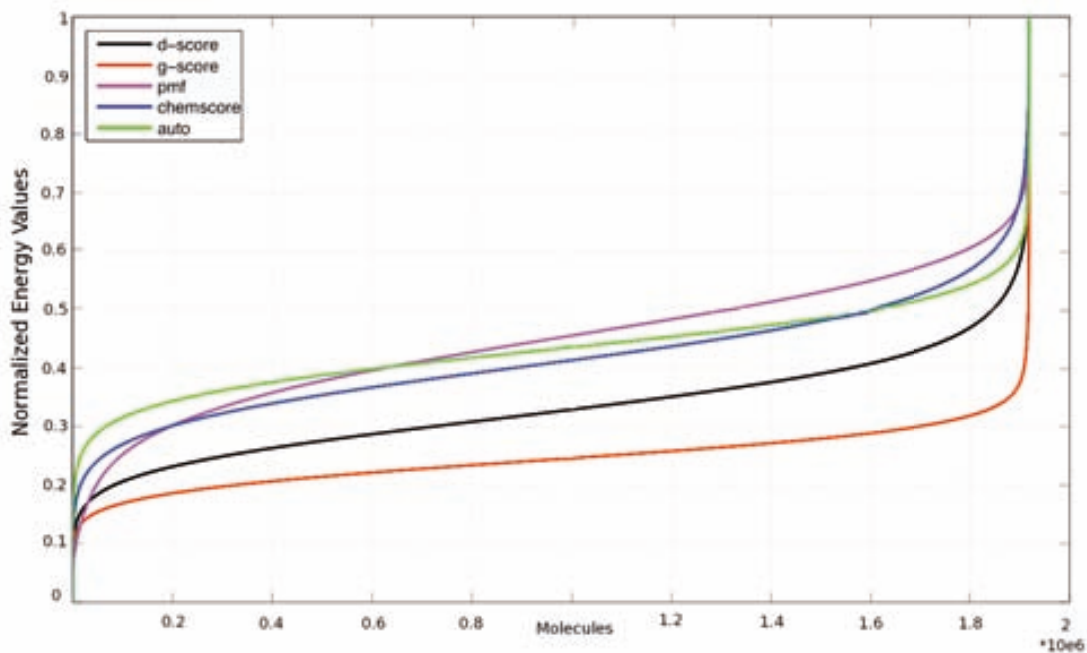
**Figure 12:** The pharmacophore filter created from three known binders of Chk1. In the figure only molecule E1 (staurosporine) is shown for clarity (excluded volumes not shown).

### 3.5.2. Docking With Autodock Vina For The Hybrid Method

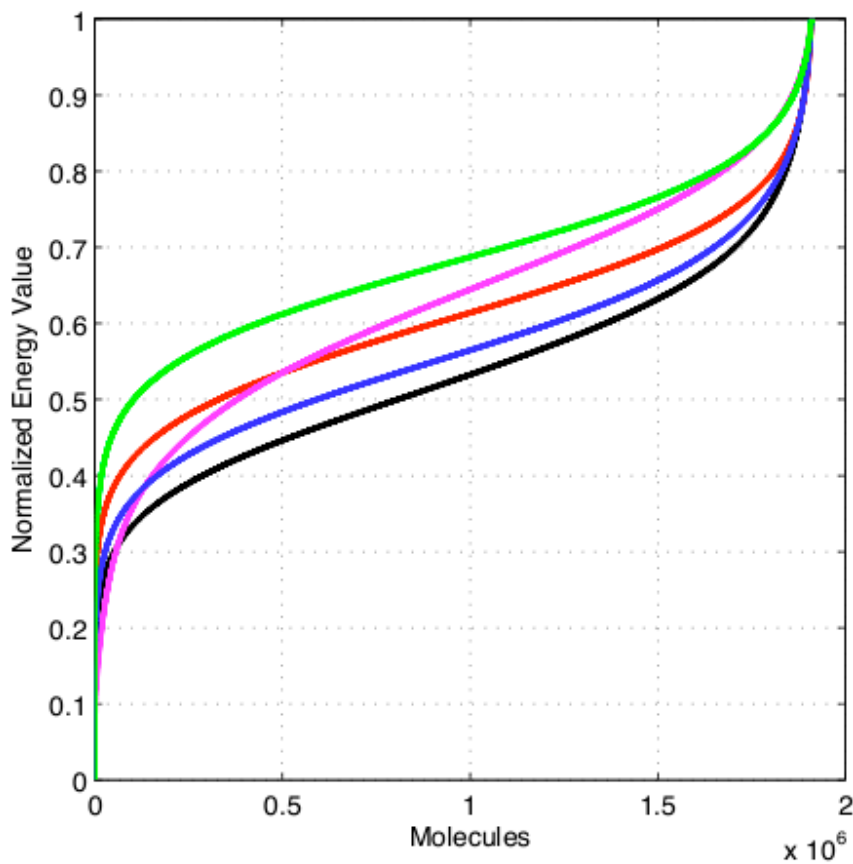
For both sets of molecules that passed the two filters, the same docking procedure was applied. The molecules that passed the pharmacophore filters were docked to the Chk1 binding site using a grid with dimensions 22 x 22 x 20 Å and 1 Å spacing. The center of the grid was decided based on the grid center used in section 3.4.3, however some small location changes were done in a way to include all pharmacophoric features. An exhaustiveness parameter of 8 was used for the docking experiments, generating 9 different poses per compound with a maximum energy difference between the best and worst displayed binding modes of 3 kcal/mol. For each compound, the pose with the lowest calculated binding free energy was saved for consensus scoring.

## 3.6. Results And Discussion

The size filter applied in the three-step docking procedure was not very successful at reducing the number of molecules for docking. This resulted in a long computation process followed by inefficient ranking. After the first docking step the normalized energy values for each scoring function show a sigmoidal pattern, as expected. However, all the plots are rather flat in the middle region because the number of molecules docked at the first docking step of the three-step docking procedure was very high, making the ranking difficult (Figure 13). As a result, a small decimal change can cause a molecule to rank much better or much worse. The size cut-off was not very strict and was chosen in relation to potential large binders, however it wasn't very effective due to the large size of the binding site. Truncating the 0.5% of the poor scoring ends from all scoring functions helped weed out the outliers, making the distinction between poorly-docked and well-docked molecules more emphasized (Figure 14).



**Figure 13:** Normalized scores after the first docking step of the three-step docking approach before 0.5 % truncation.



**Figure 14:** Normalized scores after the first docking step of the three-step docking approach after 0.5 % truncation.

Truncating the 0.5% of the poorly scoring molecules also affected the ranking of the known molecules. Table 5 shows the consensus score rank percentages of the known molecules with 0.5% truncation and without any truncation at the end of the first docking step. Even though the rankings for most of the molecules (**A1-A6**, **B1-B4** and **B6-B7**, **C1-C3** and **D1-D2**) are better without the truncation, the differences are too small to have any significant impact on the overall ranking. On the other hand 0.5% truncation boosted the rankings of molecules **E1**, **E2** and **E3** significantly: **E1** increased from being in the 17,9% to being in the 3,8%, **E2** increased from 67% to 26% and **E3** increased from 8,2% to 6,9%. Therefore, we decided to calculate consensus scores with 0.5% truncation for all the molecules in the procedure.

Molecule	Rank Percentage (%) (truncation = 0.5%)	Rank Percentage (%) (no truncation )
A1	97,5	97,3
A2	28,6	25,7
A3	82,1	79,9
A4	92,2	92,3
A5	98,8	99,4
A6	73,8	73,8
B1	88,1	87,8
B2	38,5	36,8
B3	83,7	83,6
B4	88,7	85,9
B5	77,7	77,9
B6	69,7	68,2
B7	88,3	86,7
C1	62,2	56,9
C2	37,9	35,2
C3	96,4	96,1
D1	78,5	71,6
D2	1,6	1,2
E1	3,8	17,9
E2	26	67
E3	6,9	8,2

**Table 5:** Relative ranking of known molecules with 0.5% truncation and without any truncation.

Consensus scoring was applied with the aim to improve ranking of the docked compounds. However, after 0.5 % truncation, the ranks of the known molecules according to their normalized consensus scores at the end of first step of docking were not as good as ranking only according to AutoDock energy values (Table 6). All molecules except for **A5** and **D2** ranked considerably better according to the AutoDock scoring function than the normalized consensus scoring. Even though the natural product staurosporine (**E1**) was able to make it to the second step according to the consensus scoring ranking, its rank based only on the AutoDock score was significantly better. Since only the molecules ranking in the top 60000 according to consensus scoring would follow to the second step of dockings, only **D2** and **E1** could make it in this case to the next step, out of 21 known binders. At the end of the second

step neither **D2** nor **E1** ranked in the top 1000 molecules, and therefore they couldn't proceed to the last step of docking.

Molecule	Rank at the end of step 1 (Consensus scoring rank)	Rank at the end of step 1 (AutoDock rank)
A1	1861886	1164637
A2	547653	107788
A3	1568588	1014641
A4	1760642	488160
A5	1887872	1885793
A6	1409438	860288
B1	1682831	995741
B2	736470	397051
B3	1599758	1538699
B4	1694189	941281
B5	1484617	949626
B6	1332098	745930
B7	1687037	908304
C1	1188475	228108
C2	724956	565467
C3	1841950	661120
D1	1500033	1040968
D2	31500	228107
E1	73088	641
E2	497492	27882
E3	132120	57471

**Table 6:** Ranks of known ligands at the end of the first docking step. Both consensus scoring and AutoDock scoring function ranks are the values after 0.5% truncation of the poor-scoring molecules.

The integration of pharmacophore models clearly improved the results. Both pharmacophore filters managed to reduce the number of molecules for docking quite efficiently, thus one step of dockings was enough. Even though the number of known ligands that passed the filter was low in both cases —4 and 5 for the first and second pharmacophore filters respectively—the ranks based on consensus scoring were better than the ranks from the three-step docking approach (Table 7 and Table 8).

The first pharmacophore filter let only molecules **A4**, **B4**, **D1** and **D2**, and their consensus scoring ranking after 0.5% truncation is better than their three-step docking ranking, although not as good as their ranks based on AutoDock Vina scoring (Table 7).

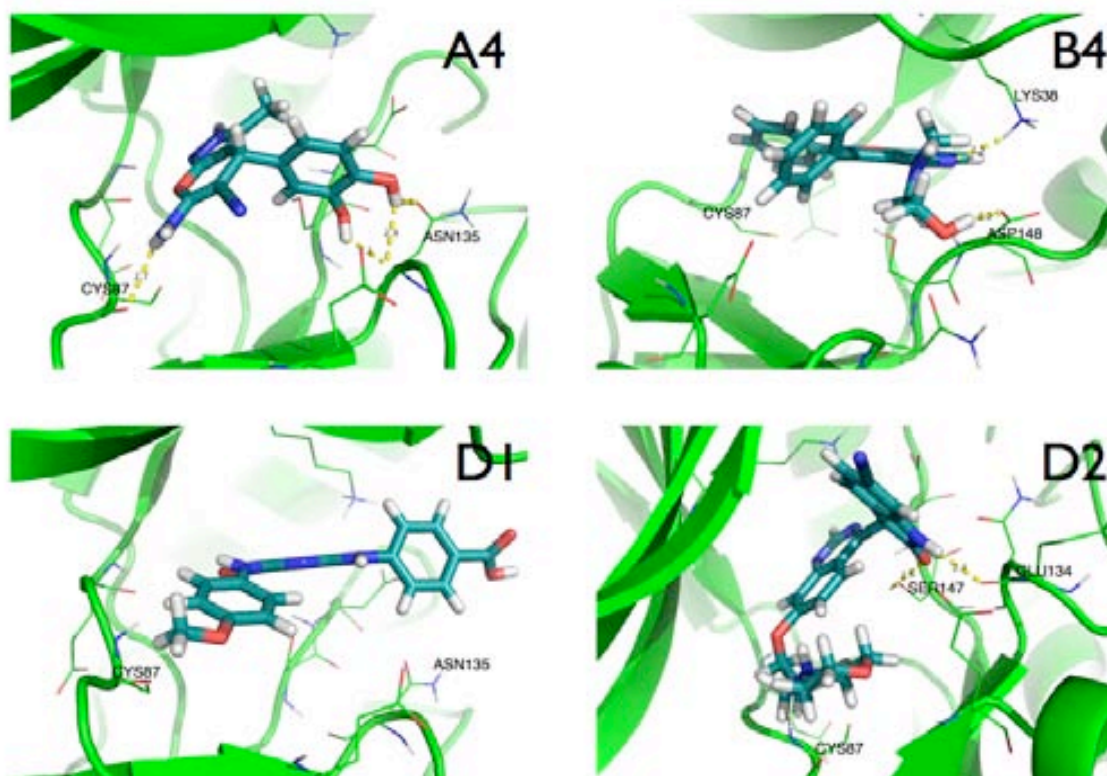
Molecule	Rank after docking (Consensus Scoring)	Rank after docking (Vina Rank)
A4	59025	25989
B4	24989	17062
D1	45880	32264
D2	4024	659

**Table 7:** Ranks of known ligands that passed the first pharmacophore filter. Both consensus scoring and AutoDock scoring function ranks are the values after 0.5% truncation of the poor-scoring molecules.

The molecules that passed the second pharmacophore filter, created from the known molecules, were **A3**, **A5**, **E1**, **E2** and **E3**. All these molecules ranked better with the hybrid approach than with the three-step docking approach. While in the case of molecules **A3** and **A5** big differences were not observed between consensus scoring rank and AutoDock Vina ranks; molecules **E1**, **E2** and **E3** ranked significantly better when only the AutoDock Vina scoring function was evaluated (Table 8). Even though the pharmacophore filter was created from these three natural molecules, the docking and scoring parts are completely independent from the pharmacophore model creation and the good docking results of **E1**, **E2** and **E3** are not direct results of the pharmacophore filter. However, the good docking results suggest an effective filtering of the small molecule library by the pharmacophore filter and an efficient differentiation of false positives from the true positives.

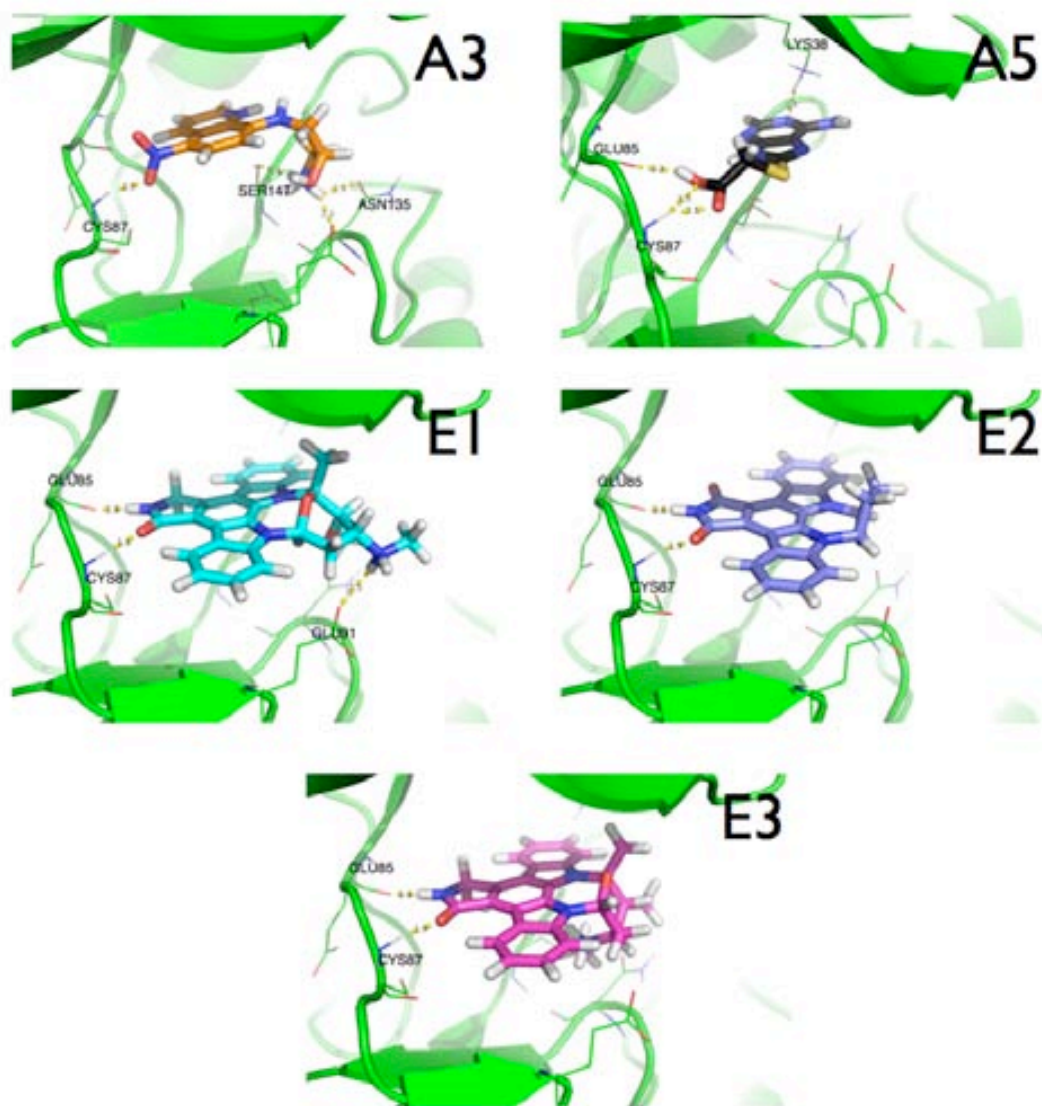
<b>Molecule</b>	<b>Rank after docking (Consensus Scoring)</b>	<b>Rank after docking (Vina Rank)</b>
<b>A3</b>	158424	161868
<b>A5</b>	170758	171777
<b>E1</b>	3346	231
<b>E2</b>	8583	7
<b>E3</b>	967	408

**Table 8:** Ranks of known ligands that passed the second pharmacophore filter. Both consensus scoring and AutoDock scoring function ranks are the values after 0.5% truncation of the poor-scoring molecules.



**Figure 15:** Docked conformations of the known ligands that passed the first pharmacophore filter.

When docked conformations of the molecules that passed the first pharmacophore filter are analyzed after docking (Figure 15), it is seen that only **A4** satisfies the requirements of the pharmacophore filter: hydrogen bonding to Cys87 and Asn135. However, both **B4** and **D2** make hydrogen bonds with the other residues: **B4** with Lys38 and Asp148, **D2** with Glu134 and Ser147. The only molecule lacking any hydrogen bonds with the target protein in the docked pose is **D1**.



**Figure 16:** Docked conformations of the known ligands that passed the second pharmacophore filter.

The molecules that passed the second pharmacophore filter make hydrogen bonds to Cys87 when docked with AutoDock Vina (Figure 16). **A3** makes one hydrogen bond to Cys87, two hydrogen bonds to Asn135 and one to Ser147. **A5** binds to Lys38, Glu85 and Cys87. Since **E3** and **E2** are derivatives of **E1**, it is not surprising that all three of them were docked in similar poses. They all make hydrogen bonds with Glu85 and Cys87, and **E1** makes an additional bond with Glu91.

### **3.7. Conclusion**

A complex problem like drug discovery needs solutions based on integrated approaches. In this study, we suggested an approach to high-throughput docking and combined it with a consensus scoring methodology we developed. In a first approach, the only criterion for the ligand-target complementarity was the size of the binding site and the ligands. However, when virtually screening large compound libraries, chemical and steric properties of the target



protein and the ligands should be taken into account and more solid properties defining ligand-target complementarity should be chosen for a directed screening. The three-step docking approach also proved that a more sophisticated and directed approach was needed, as only 2 of the 21 known ligands were ranked in positions that would be saved for the next round after the first docking step. Therefore, pharmacophore models were integrated to the methodology and the approach evolved to a hybrid method. The pharmacophore models managed to reduce the library size very efficiently, weeding out potential false positives. The hybrid approach was also successful at finding out the known ligands among almost 2000000 compounds.

## *Chapter 4*

### *The Computational Workflow*



## ***4. The Computational Workflow***

We developed a computational workflow for virtual screening and applied it to two proteins that are suggested to take part in the complex neurobiological mechanisms involved in Alzheimer's Disease<sup>128</sup>, human T-protein<sup>129</sup> and human bleomycin hydrolase<sup>130</sup>, and to acid-beta glucosidase<sup>131</sup>, involved in Gaucher's Disease pathology. The workflow includes i) pharmacophore creation from either the target protein or the known active ligand or from both, ii) pharmacophore filtering to reduce the size of the library, iii) molecular docking of the filtered compounds to the target's binding site, iv) scoring the bound poses with different scoring functions, and v) running two sets of molecular dynamics simulations on a smaller selected subset of compounds to predict their binding free energy to the target protein by the LIE method.

### ***4.1. Biological Background***

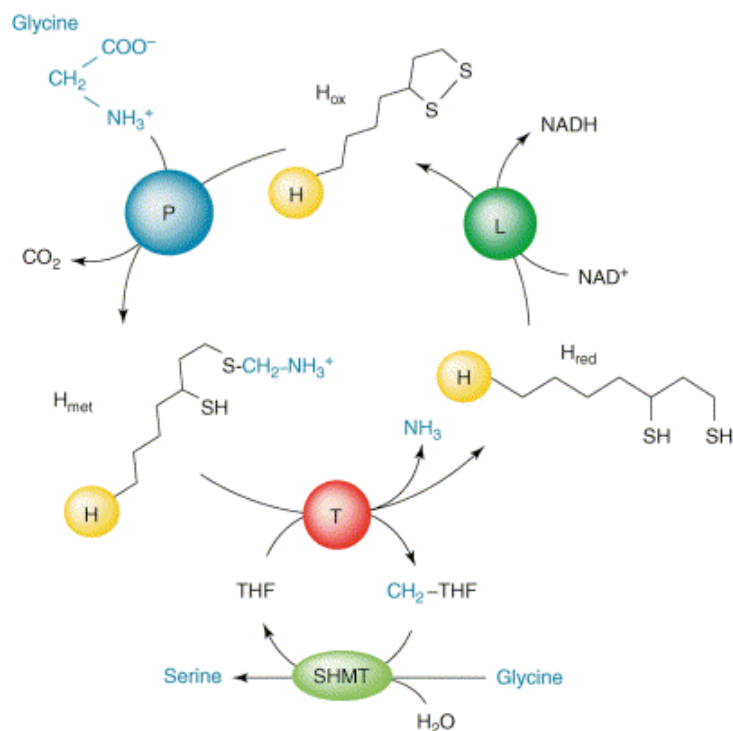
This section gives information about the structural properties of the target proteins used for this part of the study and their mode of action.

#### ***4.1.1. Human T-Protein***

Human T-protein is a component of the glycine cleavage system (GCS) and works in catalyzing the degradation of glycine.<sup>132,133</sup> A defect in any component of the GCS can abolish the overall activity of this system, resulting in elevated levels of glycine in blood and cerebrospinal fluid and leading to non-ketotic hyperglycinemia.<sup>134,135,136,137,138,139</sup>

Along with human T-protein, 3 more proteins are involved in GCS: P, L and H proteins (Figure 17). Glycine cleavage by GCS starts with P-protein catalyzing the decarboxylation of glycine, releasing CO<sub>2</sub>. As a result, the aminomethylene group binds to the H-protein. H-protein carries decarboxylated glycine to the T-protein. The T-protein has two binding sites: one for the lipoamide and one for folate. It catalyzes the transfer of a methylene carbon unit from already decarboxylated glycine attached to the lipoate cofactor of H-protein in the lipoamide binding site to the H<sub>4</sub>folate in the folate binding site, releasing ammonia and reducing the H-protein as a result. The reduced H-protein is then re-oxidized by the L-protein in a reaction that utilizes NAD<sup>+</sup> and produces NADH to complete the cycle.

The T-protein's ammonia-releasing step of this mechanism may be directly linked to steady-state levels of brain ammonia, which has led to an ammonia hypothesis of Alzheimer's Disease.<sup>140</sup> In addition, various studies showed that glycine is needed for proper functioning of NMDA receptors.<sup>141,142,143,144</sup> Changes in glycine binding to the glycine recognition sites on NMDA receptors lead to psychosis, depression and anxiety, a set of symptoms that are common to Alzheimer's disease, and inhibitors of GCS have been proposed as potential antipsychotics.<sup>145,146</sup>



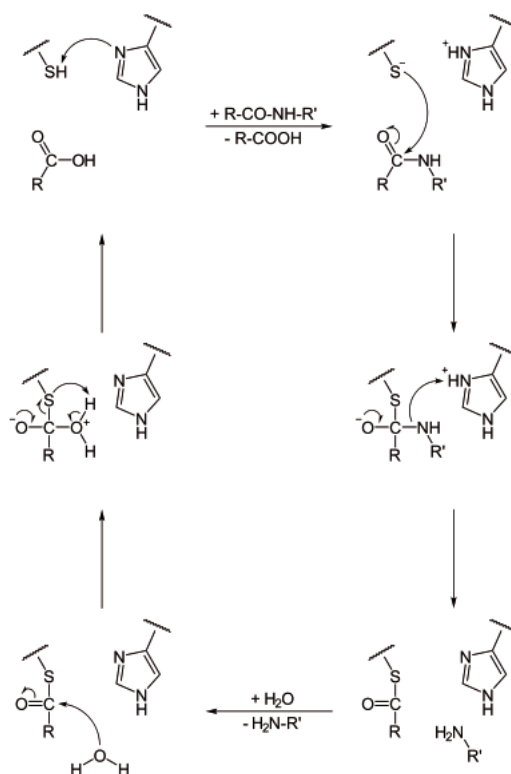
**Figure 17:** Glycine cleavage system mechanism. Human T-protein works with H, L and P proteins for degradation of glycine (Figure taken from Douce et.al.<sup>147</sup>).

#### 4.1.2. Human Bleomycin Hydrolase

Human bleomycin hydrolase is a papain superfamily cysteine protease that has the signature active site residues of this family and acts as an aminopeptidase with broad substrate specificity.<sup>148,149,150,151,152</sup> The cysteine protease mechanism of peptide bond cleavage is well defined and involves a series of steps that starts with the deprotonation of the thiol of Cys by the His side chain, which is oriented by Asn to allow this protonation. Anionic Cys sulphur makes a nucleophilic attack on the substrate carbonyl atom, releasing the amino terminus of the substrate and changing the His back to its deprotonated form. The intermediary bond linking the carboxy terminus to Cys is broken by hydrolysis to release the carboxy terminus of the substrate and to restore the enzyme back to its free form (Figure 18).

Since it deactivates the cancer therapeutic bleomycin, bleomycin hydrolase is thought to be the major cause of tumor cell resistance to bleomycin chemotherapy.<sup>153</sup> Even though bleomycin hydrolase was discovered because of its ability to detoxify bleomycin, the abundance of its homologs in different tissues of different organisms and the evolutionary conservation of the active residues propose a currently undiscovered cellular function.<sup>130</sup> Human bleomycin hydrolase has been shown to interact with human ribosomal proteins<sup>154</sup>, ubiquitin-conjugating enzyme 9<sup>155</sup> and amyloid precursor protein<sup>156</sup>. Amyloid precursor protein is the source of  $\beta$ -amyloid peptides that aggregate creating amyloid plaques of sporadic and familial cases of Alzheimer's Disease.<sup>157</sup> Various studies have indicated the presence of bleomycin hydrolase in dystrophic neurites of amyloid plaques and proven its presence in the processing of the amyloid precursor protein.<sup>158,159</sup> No specific inhibitor for Human bleomycin is known, however, the protein is inhibited irreversibly by covalent

binding of the cysteine protease specific inhibitor trans-epoxysuccinyl-L-leucylamido(4-guanidino)butane (E64)<sup>149</sup>.



**Figure 18:** Cysteine protease mechanism of peptide bond cleavage. The same mode of action is observed for human bleomycin hydrolase. All cysteine proteases have the same catalytic residues: Cys, His and Asn (Figure taken from [http://en.wikipedia.org/wiki/Cysteine\\_protease](http://en.wikipedia.org/wiki/Cysteine_protease)).

#### 4.1.3. Human Acid $\beta$ -Glucosidase

Several sporadic and genetic diseases are caused by protein misfolding.<sup>160</sup> Gaucher's Disease (GD) is a lysosomal storage disease caused by mutations in the gene (*GBA*) encoding acid  $\beta$ -glucosidase (GCase) that cause the protein not to fold into the stable form.<sup>161,162,163</sup> The disease manifests itself with symptoms like enlarged spleen and liver, liver failure, skeletal and bone disorders, anemia and in severe cases central nervous system (CNS) disorders. Mutations in GCase disrupt the degradation of glucosylceramide into glucose and ceramide resulting in accumulation of glucosylceramide in the lysosomes, causing Gaucher's Disease. Even though there are over 250 mutations related to *GBA*, the disease provoking mutations are a few prominent ones.<sup>164</sup>

GD related mutations either reduce the catalytic activity of GCase or cause a loss of protein stability during synthesis.<sup>165,166</sup> However; of all these GCase mutations, the ones that cause reduced protein stability, and therefore misfolding inside the endoplasmic reticulum, are the main reasons for GD.<sup>167</sup> Mutant GCase is mostly broken down by endoplasmic reticulum associated degradation (ERAD) due to misfolding and cannot be trafficked to the lysosome even though the remaining fractional activity of mutant GCase is still enough to hydrolyze

glucosylceramide.<sup>167</sup> It has been shown that mutant GCCase is almost as stable as the wild type in the acidic environment of lysosomes while it displays reduced stability in the neutral environment of the ER.<sup>168</sup>

There are two different types of treatment available for GD: enzyme replacement therapy<sup>169,170</sup> (ERT) with recombinant human GCCase or substrate reduction therapy<sup>171</sup> (SRT) with *N*-butyldeoxynojirimycin<sup>172</sup>. Both therapies aim to reduce the glucosylceramide stock in the lysosome, thus treating the maladies caused by its accumulation.<sup>173</sup> However, costly and life-long treatment with ERT and the abundance of side effects of SRT make it necessary to seek new therapeutic approaches.<sup>173</sup>

Small molecules that bind to misfolded proteins and guide them to correct folding by stabilizing the native state of these mutant proteins are called “pharmacological chaperones” and they have been proposed as new methods for treatment of GD and other protein-misfolding diseases.<sup>174</sup>

In the case of GCCase, these pharmacological chaperones are competitive inhibitors of the protein and bind to the mutant GCCase in the ER.<sup>175</sup> Newly synthesized GCCase is translocated into the endoplasmic reticulum, where molecular chaperones facilitate its proper folding. The molecular chaperones then dissociate from the folded enzyme, which moves to the Golgi apparatus and then to lysosomes, where the enzyme is stable and active in the acidic environment of these organelles. In the case of a mutation, misfolded GCCase molecules are degraded in the endoplasmic reticulum by ERAD. However, certain missense mutations decrease the stability of the enzyme, but the conformation of the active site is retained. Most of this type of mutant enzyme may be stabilized by pharmacological chaperones that bind to the active site of the enzyme, promote folding, and stabilize the mutant enzyme. Some of the enzyme then reaches the lysosomes, where it retains low levels of activity. In the lysosomes, the accumulated substrates displace the pharmacological chaperones and are hydrolyzed by the enzyme. Molecules like *N*-nonyl-deoxynojirimycin<sup>175</sup>, *N*-octyl- $\beta$ -valienamine<sup>176</sup>, the iminosugar isofagomine<sup>177,178</sup> and ambroxol<sup>179</sup> have been shown to increase lysosomal GCCase activity. The concept of pharmacological chaperoning makes it possible that orally administered small molecules substitute intravenous enzyme replacement therapy as the standard treatment for GD and other lysosomal storage diseases.<sup>174</sup>

## ***4.2. Methods***

This section elaborates on the details of the methods employed at different steps of the computational workflow created. Since the workflow is a cascade of steps, the output of one step is the input for the following. Here we focus on the details of the preparation of the small molecule libraries and the target proteins, pharmacophore filter creation and pharmacophore searches, docking, consensus scoring and binding free energy calculation using the LIE method.

#### 4.2.1. Preparation Of Small Molecule Libraries And The Target Proteins:

##### *The small molecule libraries:*

For this part of the work we used two small molecule libraries: VSL-1 and VSL-2. VSL-2 is an update of VSL-1 in the sense that it contains some additional vendors, and that the availability of each compound was re-checked and compounds no longer available were removed. VSL-2 was prepared in the same manner as VSL-1, as explained in section 3.3.1. VSL-1 contains 1,961,165 and VSL-2 contains 2,157,575 compounds. We used VSL-1 for screening against human T-protein, and VSL-2 for screening against human bleomycin hydrolase and GCCase.

##### *The target proteins:*

The human T-protein structure, with PDB id 1WSV<sup>129</sup>, is dimeric. However Okamura-Ikeda et.al. mentioned that dimer formation may be the result of crystal contacts. Therefore, the procedures were applied on only one chain. We also used only one chain of the homohexamer human bleomycin hydrolase (PDB id 1CB5<sup>130</sup>). The GCCase structure used for virtual screening, with PDB id 2NSX, has four chains and the competitive inhibitor isofagomine (IFG) bound.<sup>131</sup> For this study, only chain B of 2NSX was used.

All proteins were pre-processed with the Biopolymer module of SYBYL-X by Tripos<sup>34</sup> before docking: bound ligands in case of human T-protein and GCCase and water molecules were removed, hydrogens and charges (AMBER7\_F99 charge set) were added. For each target, a short minimization using the AMBER7\_F99 force field with the Powell method and termination after 500 iterations was done. These pre-processed structures were used as input for pharmacophore modeling with UNITY, docking with AutoDock Vina and consensus scoring with CSCORE<sup>87,34</sup>. However, for AutoDock Vina's file format PDBQT, the charges were replaced with Gasteiger charges<sup>125</sup> and non-polar hydrogens were merged.

#### 4.2.2. Pharmacophore Creation And Search

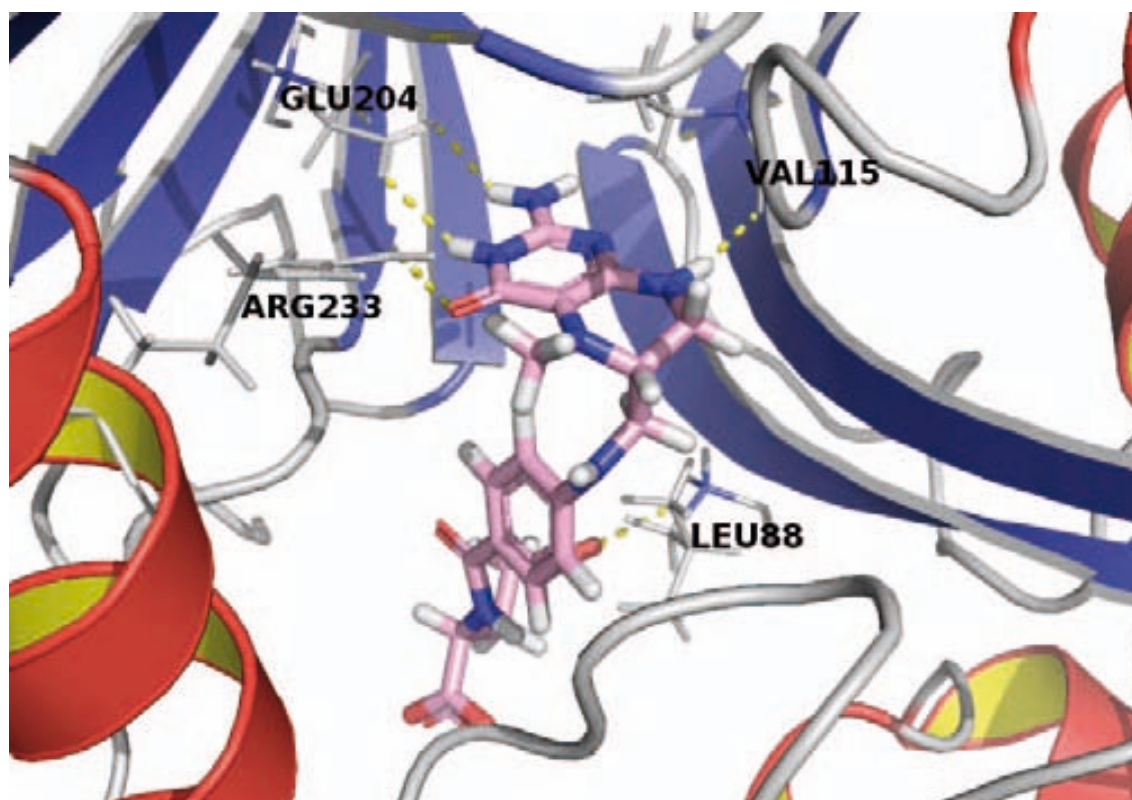
In this study, the UNITY tool from SYBYL-X was used to create pharmacophore queries and to perform 3D flexible database searches for these queries. Unlike a static 3D search, a 3D flexible search does not restrict the search to the 3D conformations as stored in the database, but generates all feasible conformations giving out only the relevant ones to the query.

##### *The human T-protein*

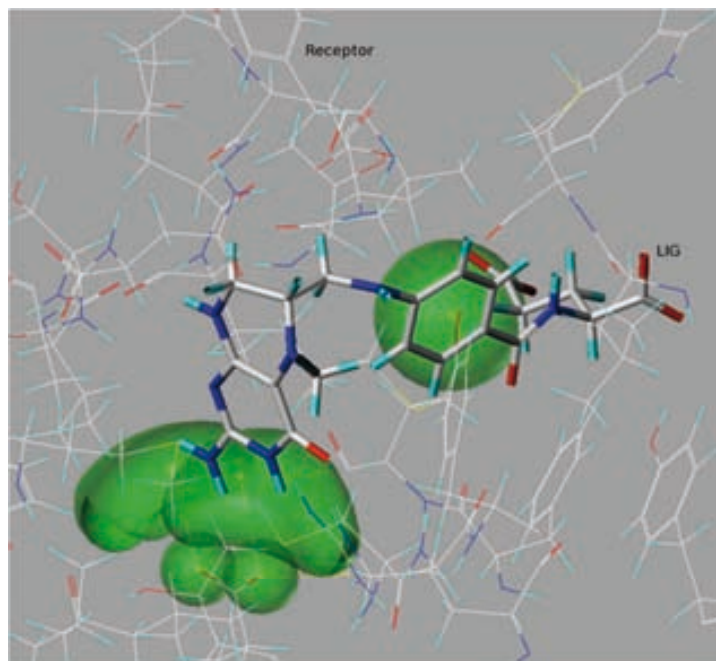
The structure of T-protein has a folate substrate analog bound in the folate binding site with an extensive network of hydrogen bonding and hydrophobic interactions, fostering an abundance of pharmacophoric features. The binding site is a hydrophobic pocket deeply buried towards the core of the protein. The residues that contribute to hydrogen bonding are Val115, Glu204 and Arg233 (Figure 19). The side chain of Glu204 makes a double hydrogen bond to the folate analog, reproducing a common behavior of folate-dependent enzymes.<sup>129</sup> Therefore, we decided to use this property of making a double hydrogen bond for creating the pharmacophore filter, defining two hydrogen bond acceptor site features on the side chain of Glu204 that require two hydrogen bond donor atom features in the candidate ligand (Figure 20). Even though the overall screening procedure is receptor-based, to make the



pharmacophore more specific the property of having an aromatic ring (van der Waals scaling factor of 2) and the directionality of hydrogen bonds were taken from the known ligand (Figure 20). Both hydrogen bond acceptor and donor features were defined with van der Waals tolerance of 1 Å. The filter was completed by adding the residues of the binding site and around as excluded volume, which represents the volume that cannot be occupied by the candidate ligands. However, the excluded volume spheres around each atom were scaled by 0.25 to allow for flexibility on the candidate ligands.



**Figure 19:** The binding site of human T-protein with the competitive inhibitor. The hydrogen bonds to Leu88, Val115, Arg233 and the double hydrogen bond to Glu204 are shown as yellow dashed lines.

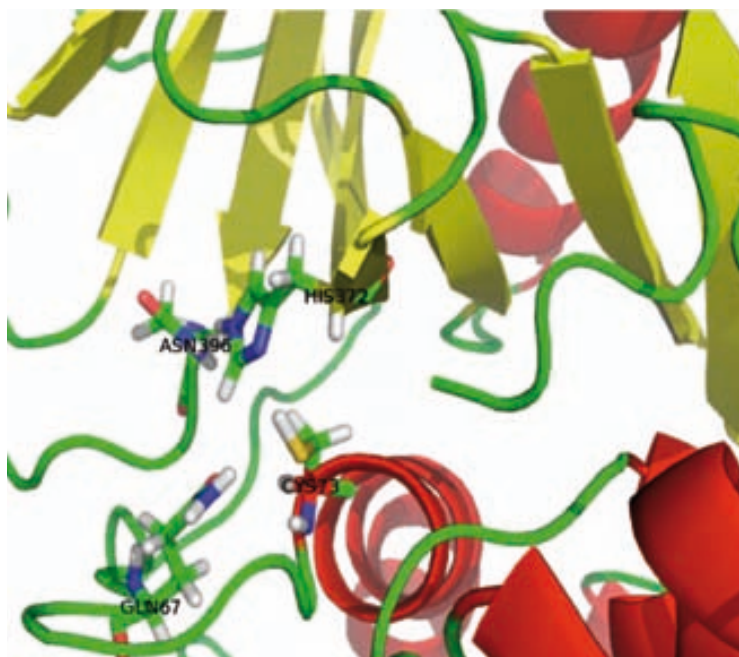


**Figure 20:** The pharmacophore filter created from the T-protein binding site residue Glu204 and from the known inhibitor. Two hydrogen bond acceptor site features on the side chain of Glu204 require two hydrogen bond donor atom features for the candidate ligand. Having an aromatic ring feature was added from the bound inhibitor (excluded volumes not shown).

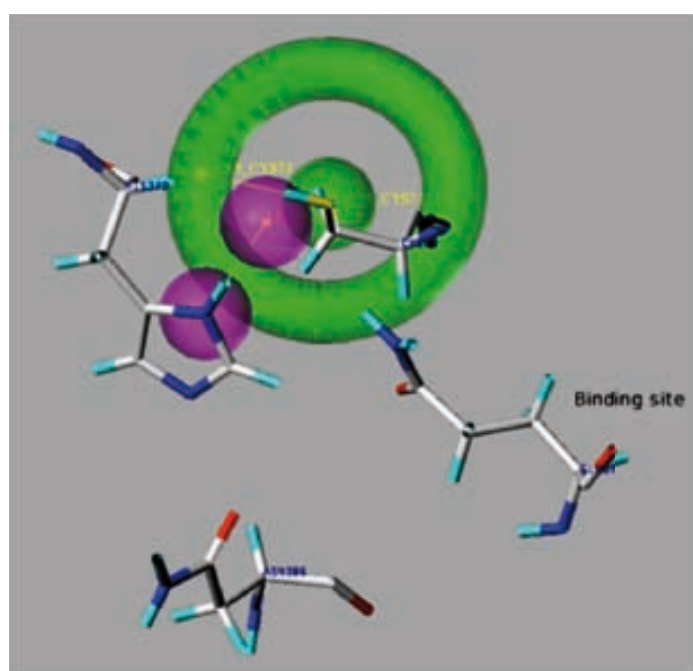
#### *The human bleomycin hydrolase*

Human bleomycin hydrolase is a cysteine protease with a ring barrel structure that has the active residues Cys73, His372 and Asn396 embedded in the central cavity (Figure 21).<sup>130</sup>

Preventing the deprotonation of Cys by His that starts the catalytic mechanism was the approach used for designing the pharmacophore filter. Therefore we defined a hydrogen bond donor site feature on Cys73 that would need a torus shaped hydrogen bond acceptor atom feature in the pharmacophore filter. Then a hydrogen bond acceptor site feature on His372 that would need a hydrogen bond acceptor atom feature in the ligand was added to the filter (Figure 22). The directionality of the hydrogen between His372 and the candidate ligand was adjusted to be towards Cys73. Finally the excluded volume feature was added from the binding site residues. The tolerance, i.e. allowed deviation from coordinates to which the features are constrained, was 1 Å for all hydrogen bond acceptor and donor features, and the van der Waals scaling was 0.25 for the excluded volume features. Although the cysteine protease inhibitor E64 makes a covalent bond with the thiol of Cys, this information wasn't used for pharmacophore design, making the process strictly structure-based.



**Figure 21:** The binding site of human bleomycin hydrolase. The catalytic residues Gln67, Cys73, His372 and Asn396 are shown.

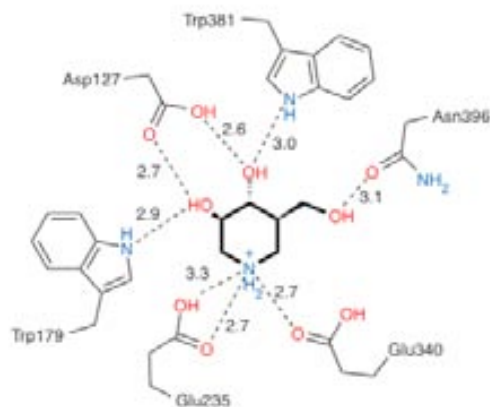


**Figure 22:** The pharmacophore filter created from the bleomycin hydrolase binding catalytic site residues Cys73 and His372 (excluded volumes not shown).

### *The human acid $\beta$ -glucosidase*

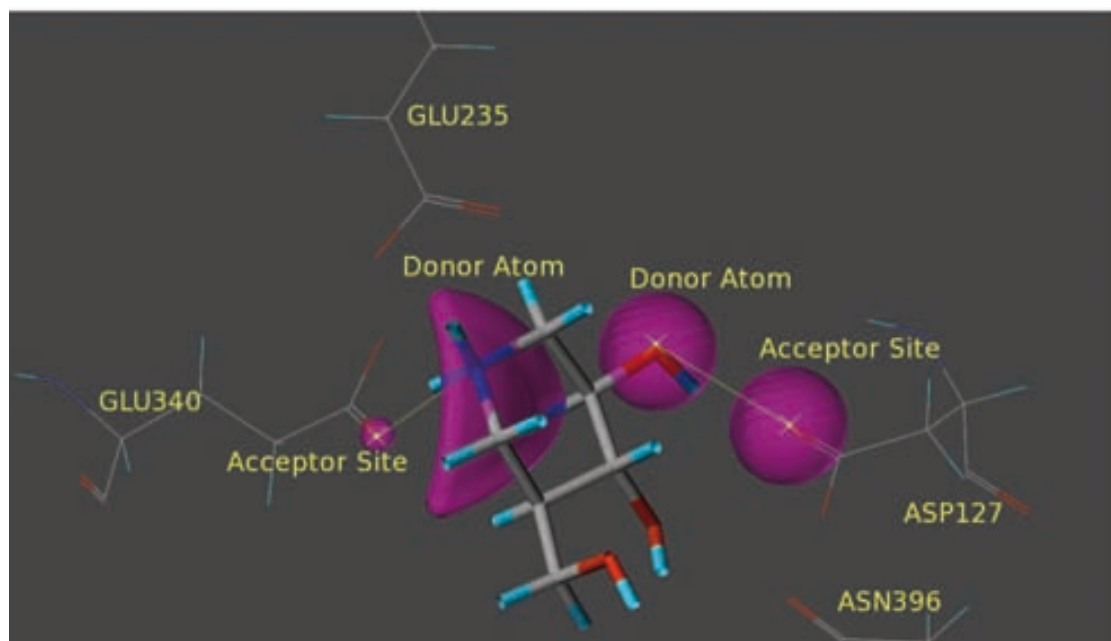
The PDB structure 2NSX has the competitive inhibitor IFG bound in the active site with an extensive network of hydrogen bonds. The imino group of IFG is stabilized by Glu235 and Glu340 while Asp127, Trp179, Trp381 and Asn396 interact with the hydroxyl groups of IFG (Figure 23).<sup>53</sup> These interactions foster an abundance of pharmacophoric features both on the

protein binding site and the known ligand (Figure 23). We created two different 3D flexible pharmacophore filters to screen the small molecule library.



**Figure 23:** Binding site of 2NSX and IFG. IFG is positioned making an extensive network of hydrogen bonds while interacting with Asp127, Trp179, Glu235, Glu340, Trp381 and Asn396. (Image taken from Farrell et. al. 53)

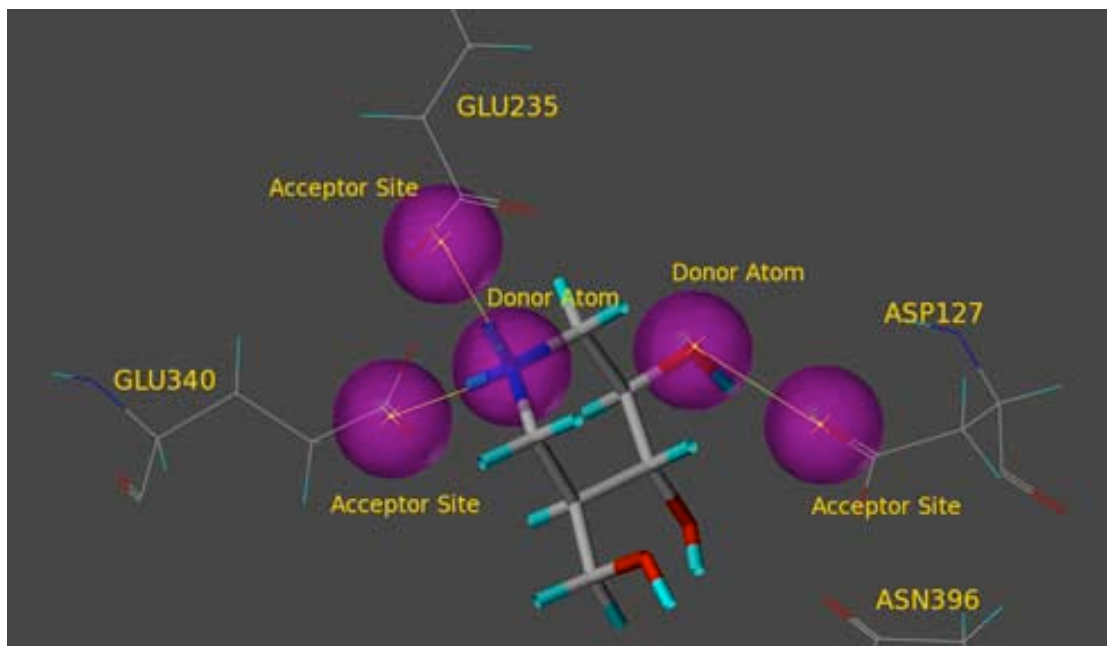
The first pharmacophore filter, **pharma1** (Figure 24) was designed from the hydrogen bond between Glu340 and the imino group of the pyranose-like ring of IFG, and from the hydrogen bond between Asp127 and the hydroxyl group of IFG. A hydrogen bond acceptor site feature (tolerance=0.3 Å) was placed on the carbonyl oxygen atom of the carboxyl group of Glu340, requiring a hydrogen bond donor atom feature (tolerance=0.3 Å) in the pharmacophore query for the candidate ligand. This feature was defined from the protein binding site. On the other hand, the pharmacophoric feature for a hydrogen bond donor atom (tolerance=1 Å) on the candidate ligand was derived from the hydroxyl group of IFG, which makes a hydrogen bond with the carboxyl oxygen of Asp127 (tolerance=1 Å). Afterwards, **pharma1** was completed by adding all binding site residues as excluded volume with the Van der Waals atom radius scaled by 0.25. Consequently, **pharma1** consists of two hydrogen bond donor features, first requiring a hydrogen bond with Glu340 and second requiring a hydrogen bond with Asp127, and excluded volume features, representing the binding site.



**Figure 24:** The first pharmacophore filter (**pharma1**) designed for GCCase. A hydrogen bond acceptor site feature was placed on Glu340, requiring a hydrogen bond donor atom feature on the candidate ligand. A hydrogen bond donor atom feature on the candidate ligand was derived from the hydroxyl group of IFG, hydrogen bonding with hydrogen bond acceptor atom feature on Asp127 (excluded volumes not shown).

The features of the second pharmacophore, **pharma2** (Figure 25) were mainly deduced from the known ligand, IFG. The hydrogen bonds between the imino group of IFG and Glu235 and Glu340, stabilizing the ring of IFG, composes the first part of **pharma2**. A hydrogen bond donor atom feature (tolerance=1 Å), making hydrogen bonds with hydrogen bond acceptor site features on Glu235 and Glu340 (tolerance=1 Å), was placed on the imino group of the pyranose-like ring. The second part of **pharma2** was derived from the hydrogen bond between a hydroxyl oxygen of IFG and carboxyl oxygen of Asp127. A hydrogen bond donor atom feature (tolerance=1 Å) requiring a hydrogen bond acceptor site feature (tolerance=1 Å) on Asp127 was located on the corresponding hydroxyl oxygen of IFG. **pharma2** was also finalized with the addition of excluded volume features as explained in **pharma1**.

While the **pharma1**-filtered library was docked to GCCase using AutoDock Vina and Surflex-Dock, only AutoDock Vina was employed for docking of the **pharma2**-filtered library.



**Figure 25:** The second pharmacophore filter (**pharma2**) designed for GCCase. A hydrogen bond donor atom feature making hydrogen bonds with hydrogen bond acceptor site features on Glu235 and Glu340 constructed the first part of the query. The second part of **pharma2** was derived from the hydrogen bond between IFG and Asp127 (excluded volumes not shown).

#### 4.2.3. Docking The Pharmacophore-Filtered Libraries To The Targets

Molecular docking was done to predict the optimum non-covalent binding of the ligands in the receptor binding sites and their corresponding binding affinities. For high throughput docking of the library to each protein target, AutoDock Vina was used. Even though AutoDock Vina allows flexibility on protein residue side chains, we left the target binding sites rigid in our study.

##### *Human T-protein and Human Bleomycin Hydrolase*

For human T-protein, the molecules that passed the pharmacophore filter were then docked to the receptor using a grid with dimensions 24 x 30 x 30 Å and 1 Å spacing. The grid was initially centered on the folate ligand. Subsequently, we manually adjusted the grid center such that it covered the whole binding cavity including its mouth, while including as little as possible of the protein surface far away from the binding pocket. Since the number of molecules that passed the pharmacophore filter for human bleomycin hydrolase is quite larger than that for the T-protein, and because the binding site of human bleomycin hydrolase is a large surface exposed area, a smaller grid was used to restrain the molecules to the binding site and to prevent molecules from being docked to irrelevant regions. The docking experiment was done in an 18 x 18 x 18 Å grid with 1 Å spacing and placed on the active site residues Cys73, His372 and Asn396. For both targets, Vina docking experiments were performed with an exhaustiveness parameter of 8, generating 9 different poses per compound with a maximum energy difference between the best and worst displayed binding modes of 3 kcal/mol. The pose with the lowest calculated binding free energy was kept for each compound for the next step of consensus scoring.

### *Human Acid $\beta$ -glucosidase*

For docking experiments with GCase, AutoDock Vina and Surflex-Dock were employed for 3 high-throughput docking experiments in total. All docking experiments were done with rigid target and flexible ligands. Docking experiment **dock1** was done with AutoDock Vina and **pharma1**-filtered library, **dock2** was done with AutoDock Vina and **pharma2**-filtered library and **dock3** was done with Surflex-Dock and **pharma1**-filtered library.

Since the binding site is a small cavity with well-defined residues, determining the center and the dimensions of the grid for the docking experiments **dock1** and **dock2** with AutoDock Vina was straightforward. The molecules were docked to the GCase structure using a grid with dimensions 20 x 16 x 16 Å and 1 Å spacing, which was placed on the bound ligand, IFG. Vina docking experiments with an exhaustiveness parameter of 8 yielded 9 different poses per compound and the pose with the lowest calculated binding free energy was kept for each compound for the next step of consensus scoring. For docking experiment **dock3** with Surflex-Dock, the computational representation of the intended binding site, protomol, was created from the binding site residues (Asp127, Trp179, Glu235, Glu340, Trp381 and Asn396) with default values. The docking experiment was done with the default parameters for the “screen” docking mode (parameters can be found in the Surflex-Dock manual<sup>181</sup>), with the only exception being the number of final poses set to 10 instead of the default value of 3. Out of 10 poses for each ligand, the pose with the best binding affinity calculated by Surflex-Dock’s scoring function was kept for consensus scoring.

#### **4.2.4. Consensus Scoring**

Consensus scoring for a given ligand-receptor complex produced by AutoDock Vina or Surflex-Dock was done as explained in section 3.3.3. First, ligand-receptor complexes output by AutoDock Vina or Surflex-Dock were rescored using D-Score, PMFScore, ChemScore and G-Score with the CSCORE program of SYBYL-X. Then, for each scoring function, all scores were normalized to values between 0 and 1 according to Equation 10:

$$S_{cut-off,F}(i) = \min\left(1, \frac{E_F(i) - E_{min,F}}{E_{cut-off,F} - E_{min,F}}\right)$$

The overall consensus score of each compound was defined by Equation 11:

$$NCS_{cut-off}(i) = \sum S_{cut-off,F}(i)$$

For the results of dockings done with AutoDock Vina (docking against human T-protein, human bleomycin hydrolase, **dock1** and **dock2** experiments for GCase),  $NCS_{cut-off}$  values of ligands vary between 0 and 5 (5 scoring functions), however the molecules docked with **dock3** experiment for GCase have  $NCS_{cut-off}$  values varying between 0 and 4 (4 scoring functions).

#### 4.2.5. LIE

Binding affinities were calculated using the LIE method.<sup>69,70</sup> This approach estimates the free energy of binding from the difference in interaction energies of the ligand with its surroundings in the protein-bound and free states. The relationship between the ligand intermolecular interaction energies and the free energy of binding is given by Equation 9:

$$\Delta G_{bind} = \alpha \Delta \langle U_{l-s}^{vdw} \rangle + \beta \Delta \langle U_{l-s}^{el} \rangle + \gamma$$

For the non-polar contribution, the coefficient has been empirically set to  $\alpha = 0.18$ . The scaling factor for the polar contribution was initially derived from the linear response approximation ( $\beta = 0.5$ ) but has subsequently been found, from free energy perturbation (FEP) calculations, to depend on the chemical nature of the ligand.<sup>70</sup> According to that classification, ligands with net charge maintain the value  $\beta = 0.5$  associated to the linear response approximation, while uncharged ligands have an empirical associated value of  $\beta = 0.43$ ,  $\beta = 0.37$  or  $\beta = 0.33$ , depending on whether they have 0, 1 or 2 or more hydroxyl groups, respectively. We have followed this classification to choose the appropriate  $\beta$  parameter for each ligand in the database. To do so, we determined the numbers of positive and negative charges and the number of hydroxyl groups of each compound in our databases using the `dbcompute` program of UNITY with `-calctype patterncount`. Finally,  $\gamma$  is a constant term obtained by fitting experimental data to LIE data with the purpose of fixing the scale for absolute binding free energies. The nature of this parameter has been related to several descriptors of the binding site, such as its hydrophobic nature.<sup>182</sup> Since we were mainly interested in prioritizing compounds in our libraries for experimental testing we were primarily interested in relative binding free energies, and thus  $\gamma$  was set to 0.

As stated below, the energies that enter into the LIE equation are averaged interaction energies of the ligand with its surroundings obtained from separate MD simulations of the ligand in water or bound to the solvated protein system (with initial coordinates of protein and ligand obtained from docking experiments). All MD simulations have been performed with the program Q<sup>183</sup> and the OPLS force field implemented therein.<sup>184</sup> Since many of the parameters as well as topologies needed for the ligands were not present in the original version of the force field, an automated parameterization protocol was followed. First, the MOL2 files used for CSCORE were converted to Schrödinger's Maestro format (.mae) using the `mol2convert` utility of the Schrödinger suite. Next, each ligand was energy-minimized with the `bmin` program of Macromodel<sup>185</sup> and the OPLS parameters and topology information generated by that program were translated into the syntax required by the program Q, using a set of *ad hoc* scripts.

MD simulations in Q were performed using spherical boundary conditions, thus a definition of a solvation sphere around the ligand was required. In our two cases, the centers of the spheres were determined manually, making small adjustments to the docking grid centers used. The sphere radii were calculated according to the diameter of the largest compound as docked into the protein binding site, according to Equation 12:

**Equation 12**

$$r_s = \left[ \frac{1}{2} L_{\max(\text{ligand})} + c \right]$$



where  $r_s$  denotes the radius of the sphere,  $L_{max(ligand)}$  is the diameter of the largest docked compound and  $c$  is a constant, which was set to 14 Å for this study. This ensures a margin of at least 14 Å of explicit environment around every atom of a ligand centered in the sphere. The same size of the sphere was used for the protein-bound and the protein-free simulations. The surface of this sphere was subjected to radial and polarization restraints<sup>186</sup> in order to mimic bulk water at the sphere boundary. Non-bonded interaction energies were calculated up to a 10 Å cutoff, except for the ligand atoms, for which no cut-off was used. Beyond the cut-off, long-range electrostatics was treated with the local reaction field (LRF) multipole expansion method.<sup>187</sup> Protein atoms outside the simulation sphere were restrained to their initial positions, and only interacted with the system through bonds, angles and torsions. The ionization states of titratable residues inside the simulation sphere were manually assessed, in order to obtain neutral simulation systems in the protein simulation, which is needed to compare the ligand-surrounding energies between bound and free states. Any titratable residues closer than 3-5 Å to the boundary of the solvation sphere, as well as those outside the solvent sphere, were modeled as neutral because of the lack of dielectric screening. Even though only movement within the sphere is allowed, amino acids with all their atoms further than  $r_s + 2$  Å away from the sphere center were removed.

For both the protein-bound and the protein-free simulations, an initial heating and equilibration MD simulation was carried out before the data collection phase, starting with a very short time step of 0.1 fs, a strong coupling to a temperature bath of 1 K and positional restraints of 25 kcal/(mol·Å<sup>2</sup>) on all non-hydrogen protein and ligand atoms in the case of the protein-ligand complex simulation. The system was then gradually heated up to 310 K during 95.5 and 65.25 ps for the protein-bound and the protein-free simulations, respectively, in which the bath coupling was relaxed to a final value of 100 fs, the timestep was increased to 1 fs and the force constant of the positional restraints was gradually lowered to 0. Detailed parameters can be found in Table 9 and Table 10.

A production-run molecular dynamics simulation then followed for 500 ps at 310 K (100 fs coupling time) with a time step of 1 fs. In the case of the protein-free simulation, the center of the ligand was restrained to the center of the solvation sphere with a force constant of 5 kcal mol<sup>-1</sup> Å<sup>-2</sup>. Energies were collected at regular intervals of 25 fs. Energy averaging was performed on this collection period, and stability was addressed by an estimation of the convergence errors of the potential energies of the ligand with its surroundings.

For the estimation of the binding free energies according to Equation 9, only the period 100-500 ps of the production-run simulations were considered. Convergence of the simulations was assessed by calculating the difference between the LIE energy calculated over the periods 100-300 ps and 300-500 ps of the production-run simulations. We refer to this difference as the LIE error.

Free State	Number of steps	Step size in fs	T (K)	Coupling (fs)	SHAKE for solvent	Force constant of restraint on heavy atoms (kcal/mol Å <sup>2</sup> )
0 – 0.25 ps	2500	0.1	1	0.1	off	n.a.
0.25 – 2.75 ps	2500	1.0	50	5	on	n.a.
2.75 – 5.25 ps	2500	1.0	150	50	on	n.a.
5.25 – 15.25 ps	10000	1.0	310	50	on	n.a.
15.25 – 65.25 ps	50000	1.0	310	100	on	n.a.

**Table 9:** Parameters used for free state simulation of the ligand in a sphere filled with water in the initial equilibration period.

Bound State	Number of steps	Step size in fs	T (K)	Coupling (fs)	SHAKE for solvent	Force constant of restraint on heavy atoms (kcal/mol Å <sup>2</sup> )
0 – 0.5 ps	5000	0.1	1	0.1	off	25
0.5 – 5.5 ps	5000	1.0	50	5	on	10
5.5 – 15.5 ps	10000	1.0	150	5	on	5
15.5 – 25.5 ps	10000	1.0	310	20	on	2
25.5 – 45.5 ps	20000	1.0	310	100	on	1
45.5 – 95.5 ps	50000	1.0	310	100	on	0

**Table 10:** Parameters used for bound state simulation of the ligand in the protein binding site in the initial equilibration period.

### 4.3. Results And Discussion

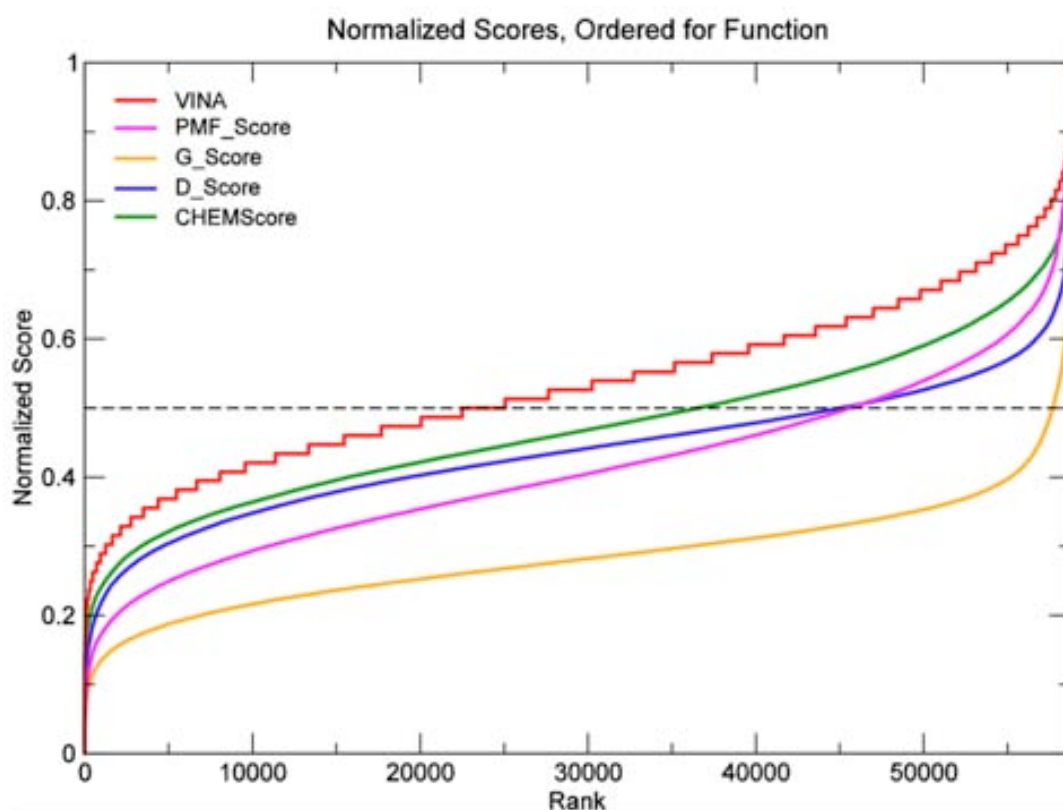
#### 4.3.1. Results For Human T-Protein And Human Bleomycin Hydrolase

The 3D flexible search with the pharmacophore filter designed for human T-protein reduced the size of VSL-1 from 1961165 to 58699 molecules. Even though the pharmacophore filter was quite tolerant, the narrow tunnel-like hydrophobic binding site and the use of pharmacophoric features from the known ligand allowed a library reduction by a factor of 30. However, in the case of human bleomycin hydrolase library reduction was not as efficient. The binding site of bleomycin hydrolase is part of a large solvent-exposed cavity. This may cause human bleomycin hydrolase to have little substrate specificity, similar to its yeast homolog.<sup>130</sup> Therefore, the pharmacophore filter created for bleomycin hydrolase was rather tolerant with fewer excluded volume features. As a result, 471198 out of the 2157575 molecules of VSL-2 passed the filter. The solvent-exposed binding site of bleomycin hydrolase also brought extra challenges for determining the parameters for the grid size used for docking. Using a big grid covering the whole cavity could have caused docking to other regions than the binding site region. Therefore, the grid chosen was as big as possible to cover the binding site residues and the docked ligands while preventing irrelevant binding.

After the docking calculations for both target proteins were completed, we selected the binding pose with the lowest energy calculated by AutoDock Vina for each ligand. These poses were subsequently rescored with the four additional scoring functions implemented in CSCORE, and the scores were normalized.

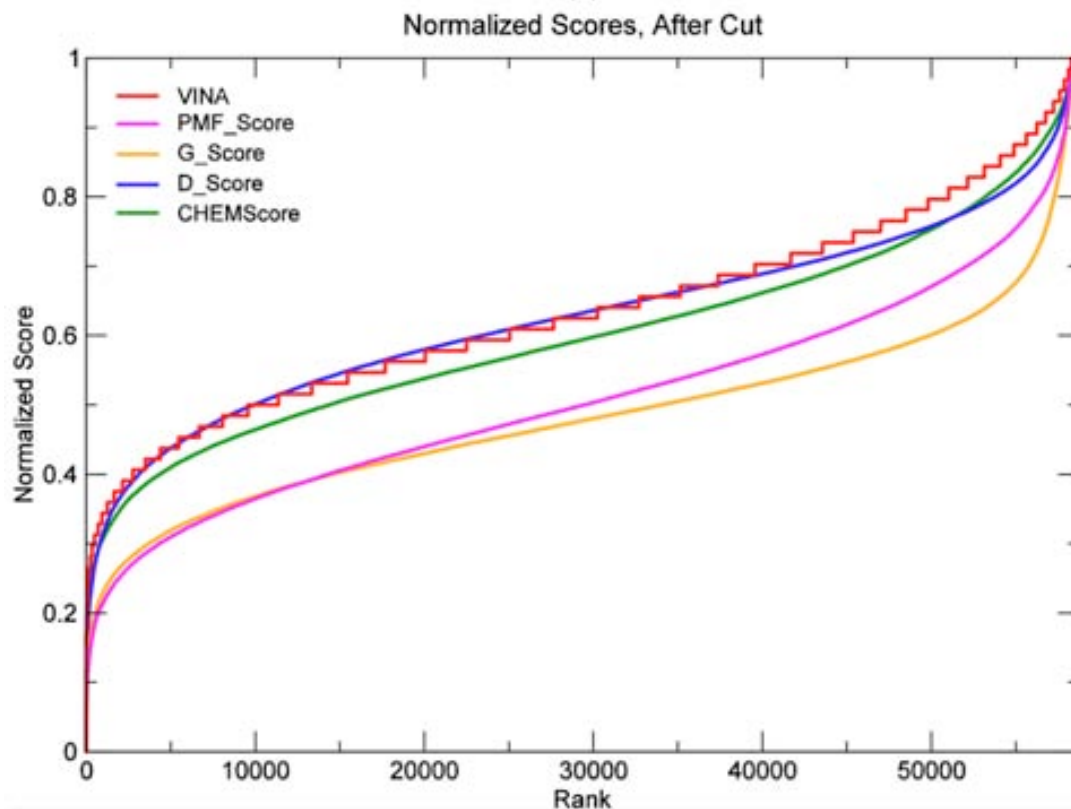
In the case of human T-protein, normalizing the scores produced by CSCORE and AutoDock Vina to values between 0 and 1, with a truncation cut-off of 100% in the normalization procedure, produced similarly shaped sigmoidal curves for all scoring functions. However, the individual scoring functions showed trails of high values to largely different extents, most pronounced in the case of G-Score (Figure 26). These values at the poor-scoring end are problematic when trying to combine the different scores to reach a consensus. Using a rank-by-vote strategy for consensus scoring based on scores being within the top  $n\%$  of the obtained score range, with a frequently-used “vote cut-off” of 0.5, results in G-Score voting

for more than 99% of the molecules, and PMFScore and D-Score for around 75% of the molecules, rendering these scoring functions nearly redundant. In turn, with a smaller value for the vote cut-off, Vina, ChemScore and D-Score vote for only a few percent of the molecules, rendering them too decisive. Using the sum of normalized scores to obtain a consensus is also problematic in this case, because the different slopes of the curves in the range of intermediate ranks results in largely different decisive power of the individual methods. For example, compared with a compound ranking around 1000, a compound ranking 20000 receives twice as large a “score penalty” with Vina as it does with G-Score. As a result, Vina has a much larger decisive power than G-Score.



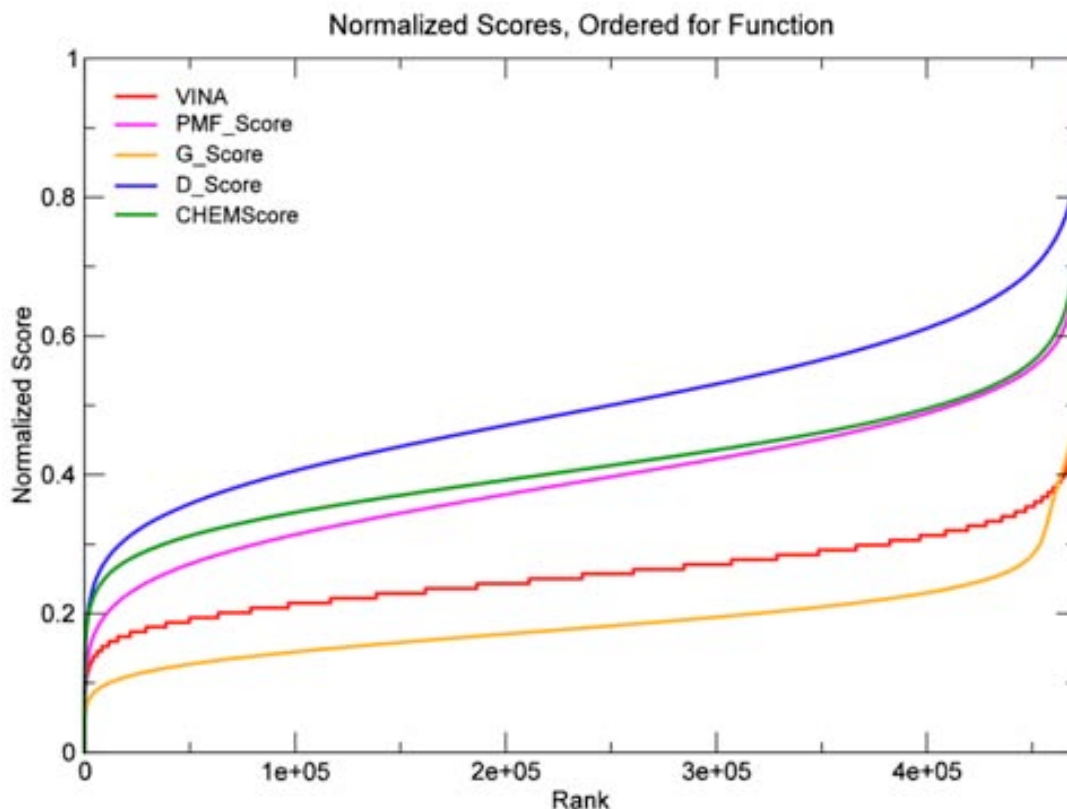
**Figure 26:** Normalized scores calculated without truncation with the five scoring functions Vina, PMFScore, G-Score, D-Score and ChemScore for human T-protein. The scores are plotted against compound ranks obtained with the individual scoring functions. The dashed line marks the frequently used vote cut-off of 0.5.

Therefore, an adjustment excluding poorly scoring compounds for each energy function was done with our normalization procedure using a truncation cut-off of 99.5%. As Figure 27 shows, the discrepancy between scoring functions has become less pronounced after this adjustment. In addition, the slopes of the curves for intermediate ranks have increased, improving the distinction between well-scored and poorly scored molecules. Hence, the sum of normalized scores calculated with the 99.5% truncated normalization procedure was used as the final normalized consensus score, *NCS*, for each compound. *NCS* values ranges between 0 and 5, the values closer to 0 representing well-scored compounds and the values closer to 5 representing poorly scored compounds.



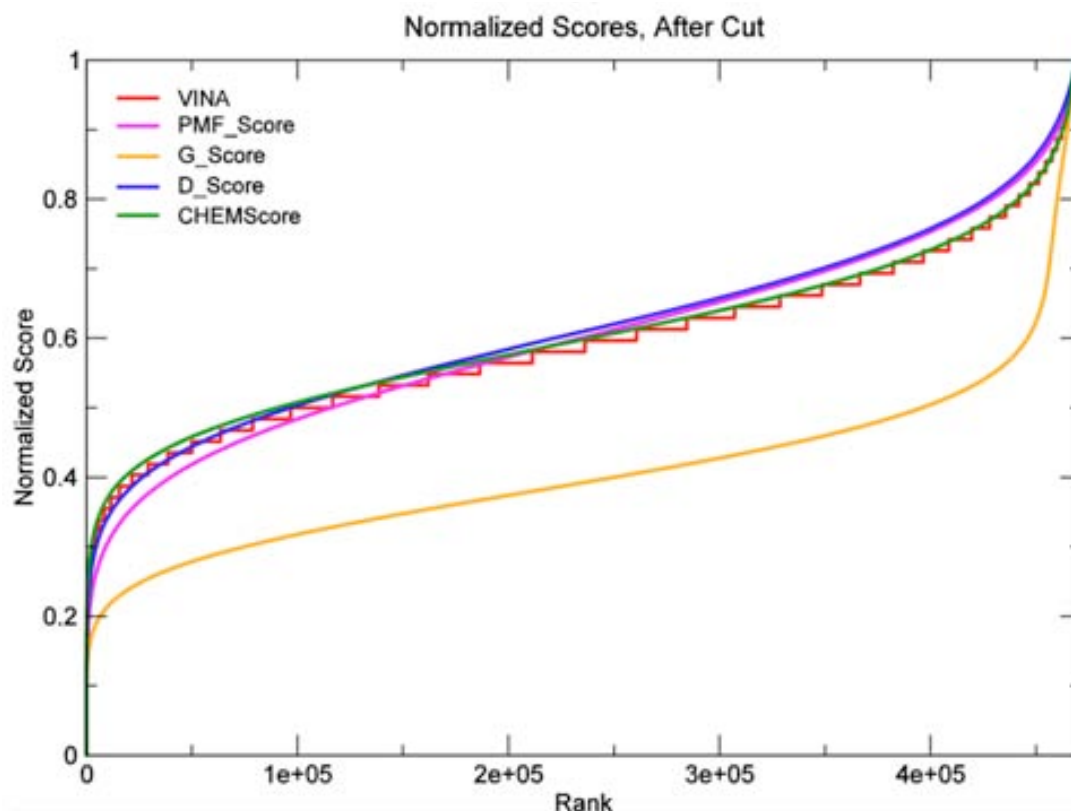
**Figure 27:** Normalized scores calculated with the second normalization procedure against compound rank (after the 0.5% worst-scoring compounds were excluded) for human T-protein.

For human bleomycin hydrolase, the normalization procedure without truncation (truncation cut-off=100%) produced results that show a similar pattern as obtained for human T-protein, however with more pronounced differences between individual scoring functions. Figure 28 shows that the ability of the normalized AutoDock Vina score and G-Score to distinguish “well docked” molecules from “poorly docked” ones is strongly affected by few compounds receiving very high scores.



**Figure 28:** Normalized scores calculated without truncation with the five scoring functions Vina, PMFScore, G-Score, D-Score and ChemScore for human bleomycin hydrolase.

To bring the scoring functions closer to each other, the normalization was done with the truncation cut-off set to 99.5%. The resulting distributions of normalized scores after the truncation were almost the same for all scoring functions except G-Score. The reason for this behavior might be that G-Score scores poorly docked molecules with a very high penalty thus losing sensitivity for fairly or slightly well docked molecules. However, making G-Score to converge to other functions would need excluding at least 5% of the molecules (around 20000 molecules). Therefore, it was decided to stop at 0.5 % truncation and calculate the normalized consensus scores at this point (Figure 29).



**Figure 29:** Normalized scores calculated with the second normalization procedure against compound rank (after the 0.5% worst-scoring compounds were excluded) for human bleomycin hydrolase.

Once the normalized consensus scores were calculated for each compound the compounds were ranked accordingly. For both target proteins, the top-ranking 5000 compounds were chosen for the next step of binding free energy estimation through LIE simulations.

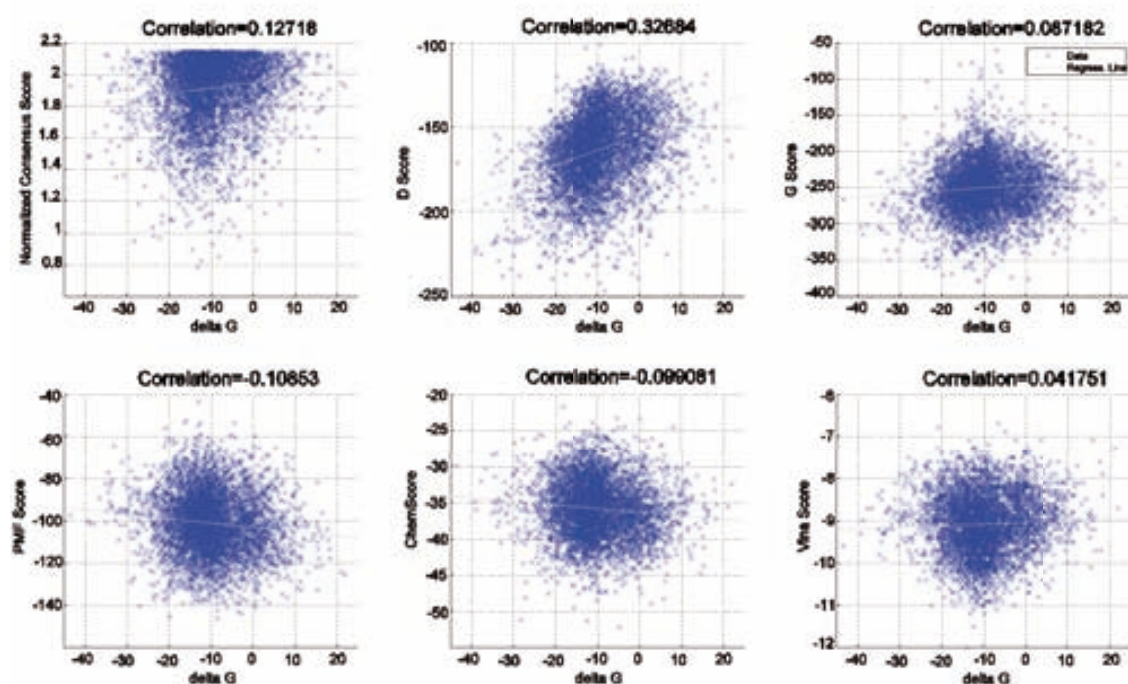
This step of LIE simulations mainly consists of five tasks: deciding the size and location of the solvation sphere, defining the protonation states of ionizable amino acids, deriving force field parameters and topologies for the compounds for the molecular dynamics simulations, defining the simulation parameters, and finally performing the simulations and binding free energy estimation for each compound-protein pair. The size of the solvation sphere depends exclusively on the size of the compounds (see section 4.2.5). For human T-protein and human bleomycin, the largest diameters of any of the 5000 selected compounds in their docked conformations were 24.27 Å and 23.1 Å, respectively, thus defining the radius of the solvation sphere as 27 Å and 26 Å, respectively, according to Equation 12. In contrast to the bleomycin hydrolase system, with a binding site that is not buried, the solvation sphere enclosed a big part of human T-protein. A very important and labor-intensive step is the manual adjustment on the protonation states of the protein's ionizable amino acids. The electrostatic environment in both sets of molecular dynamics simulations—protein-bound and protein-free—must be as similar as possible for an accurate comparison. Because the protein-free simulations take place in a sphere filled with water thus carrying no net charge, the inside of the sphere located on the protein binding site should also be electrostatically neutral. Starting from a state with all ionizable residues neutralized, we tried to ionize as many residues as possible inside the sphere and at least 5 Å from the sphere surface, prioritizing residues close to the binding site surface, while keeping the entire system electrostatically

neutral. Pairs of residues involved in salt bridges were only ionized or kept neutral together. Thus, for the human T-protein system the following titratable residues inside the simulation sphere were charged: Lys123, Lys266, Lys363, Arg6, Arg45, Arg66, Arg82, Arg191, Arg194, Arg233, Arg237, Arg267, Arg268, Arg269, Arg290, Arg291, Arg292, Asp52, Asp77, Asp100, Asp101, Asp124, Asp201, Asp234, Asp248, Asp250, Glu71, Glu80, Glu109, Glu200, Glu204, Glu239, Glu251 and Glu257. In the human bleomycin hydrolase system the charged residues were: Lys68, Lys107, Lys162, Lys309, Lys330, Lys335, Lys359, Lys405, Arg72, Arg83, Arg110, Arg175, Arg176, Arg362, Arg393, Asp35, Asp38, Asp106, Asp143, Asp179, Asp327, Asp401, Glu60, Glu96, Glu167, Glu172, Glu367, Glu395, Glu400 and Glu421.

Although 5000 molecules had been chosen for molecular dynamics simulations with human T-protein, a small number of the molecules failed during different stages of the simulation. In the end, binding free energies for 4983 molecules were calculated with the LIE method and the correlation analyses between scoring functions were done using Pearson's correlation coefficient. The comparison between LIE results and different scoring functions for human T-protein show that D-Score gave the best correlation with LIE results; however, the correlation was not very significant (Figure 30). Neither the remaining scoring functions nor the normalized consensus score showed any significant correlation with binding free energies predicted with LIE. Table 11 shows the correlations between the normalized consensus score and the scoring functions, and also the correlations of different scoring functions with each other. The only significant correlation was between D-Score and G-Score among individual scoring functions.

	NCS	D-Score	G-Score	PMFScore	ChemScore	VinaScore
NCS	1	0.53	0.52	0.29	0.49	0.26
D-Score	0.53	1	0.55	-0.14	0.04	-0.32
G-Score	0.52	0.55	1	-0.24	0.16	-0.24
PMFScore	0.29	-0.14	-0.24	1	-0.14	0
ChemScore	0.49	0.04	0.16	-0.14	1	0.05
VinaScore	0.26	-0.32	-0.24	0	0.05	1

**Table 11:** Pearson's Correlations of scoring functions with each other and normalized consensus score of the selected molecules for human T-protein.



**Figure 30:** Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for human T-protein.

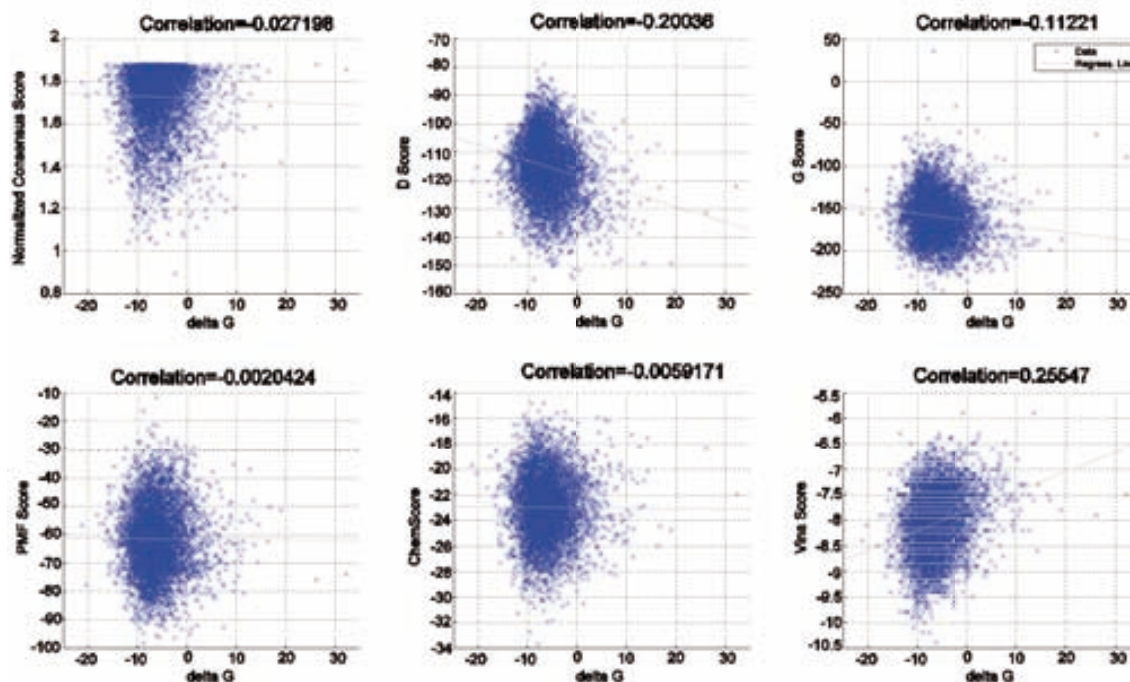
Also in the case of human bleomycin hydrolase, some ligands of the selected 5000 failed during the simulations. Binding free energies of 4960 molecules could be calculated in the end. Even though the normalized scores were correlated with individual scoring functions, binding free energies calculated with LIE only showed minor correlation with the AutoDock Vina score (Figure 31).

In both examples, there is hardly significant correlation between different scoring functions (Table 11 and Table 12). However, this is not unexpected taking into account the approximations applied by scoring functions. The other reason for the lack of correlation lies in the different characteristics of the scoring functions used. The scoring functions of AutoDock Vina and ChemScore are empirical and PMF Score is a knowledge-based scoring function, which means that they were trained against a set of protein-ligand complexes and their performances are distinctly dependent on these training sets. On the other hand, it is known that force-field based scoring functions like D-Score and G-Score usually overestimate the binding affinities.

	NCS	D-Score	G-Score	PMFScore	ChemScore	VinaScore
NCS	1	0.32	0.31	0.25	0.4	0.28
D-Score	0.32	1	0.31	-0.31	-0.19	-0.24
G-Score	0.31	0.31	1	-0.28	-0.17	-0.11
PMFScore	0.2	-0.31	-0.28	1	-0.03	-0.26
ChemScore	0.4	-0.19	0.17	-0.03	1	0.05
VinaScore	0.28	-0.24	-0.11	-0.26	0.07	1

**Table 12:** Pearson's Correlations of scoring functions with each other and normalized consensus score of the selected molecules for human bleomycin hydrolase.

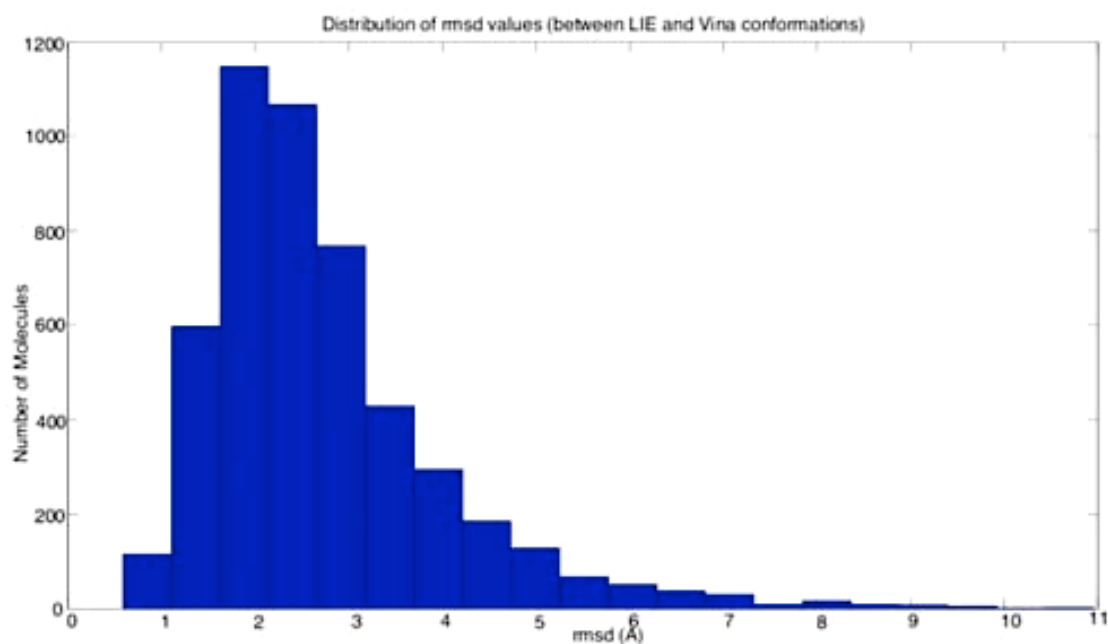




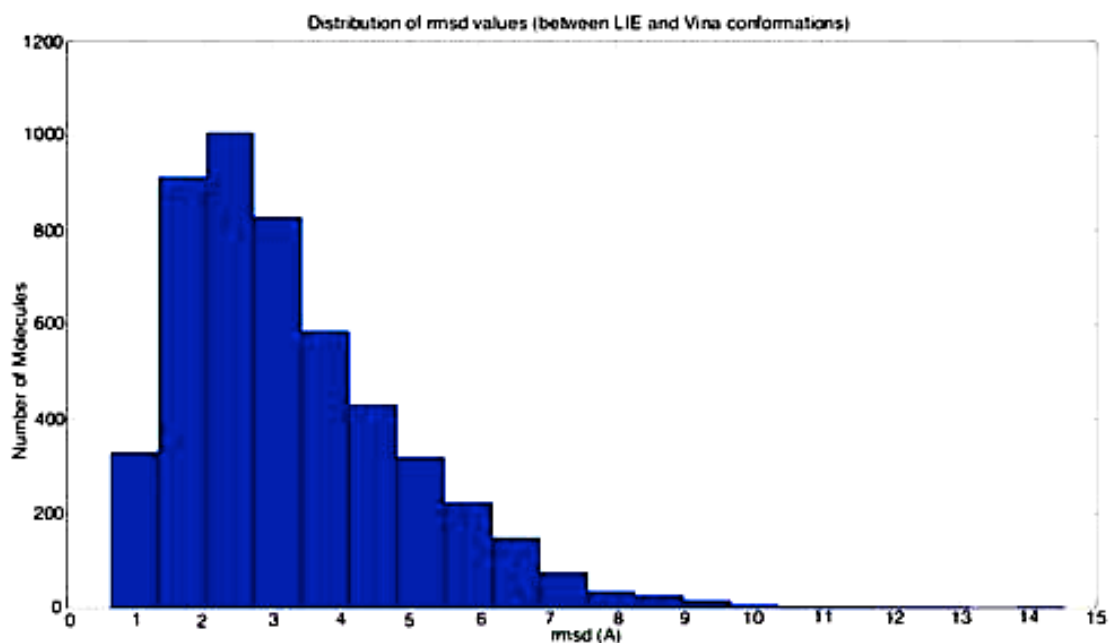
**Figure 31:** Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for human bleomycin hydrolase.

We also compared the binding modes output by the docking and the molecular dynamics simulations done in the solvated spherical system centered on the protein binding site. In the dockings, ligand positions were restricted to the grid in the protein binding site. Even though the grids were large enough for the ligands to move freely, they were not as spacious as the water filled sphere used for LIE simulations. Ligands couldn't get out of the allowed grid during docking. However, in LIE simulations, the solvation sphere did not only cover the binding site residues but also all the residues that were at most 27 and 26 Å away from the sphere center for the T-protein and bleomycin hydrolase, respectively, enabling a larger area for the ligands to explore.

For the human T-protein, the root-mean-square difference (rmsd) values between the docked conformations and final simulation conformations are mostly distributed between 1 and 4 Å (Figure 32). Around 90% of 4960 ligands fall into this area. This means that binding modes created by the docking are fairly similar to those output by the simulations. In the case of the bleomycin hydrolase, though, the percentage of molecules with similar docking and simulation binding modes is lower; around 62% of the rmsd values are in the range of 1 to 4 Å (Figure 33). This might be caused by the large solvent-exposed binding site of bleomycin hydrolase. The binding site of the T-protein is narrow and buried, restricting the molecules to bind in a certain way. The other factor for larger displacement values for ligands binding to the bleomycin hydrolase may be the existence of water molecules. While water molecules cannot fill the narrow binding site of the T-protein along with the ligand, the solvent-exposed binding site of the bleomycin hydrolase can accommodate water molecules in competing interactions.



**Figure 32:** Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) for human T-protein.

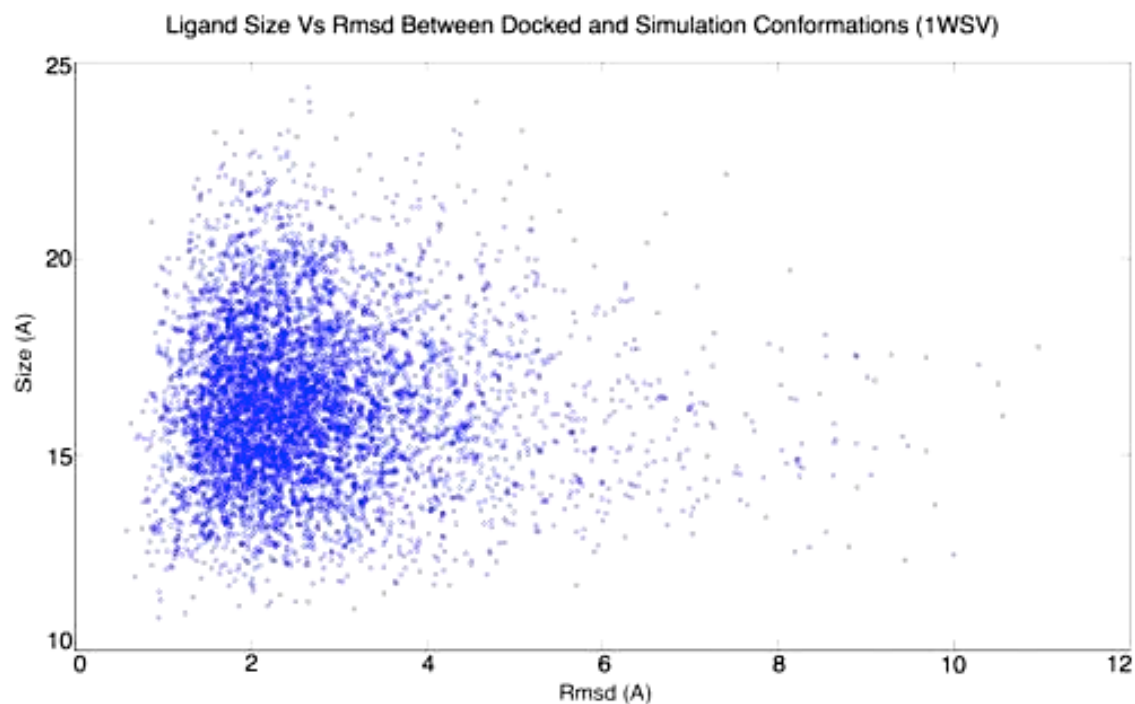


**Figure 33:** Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) for human bleomycin hydrolase.

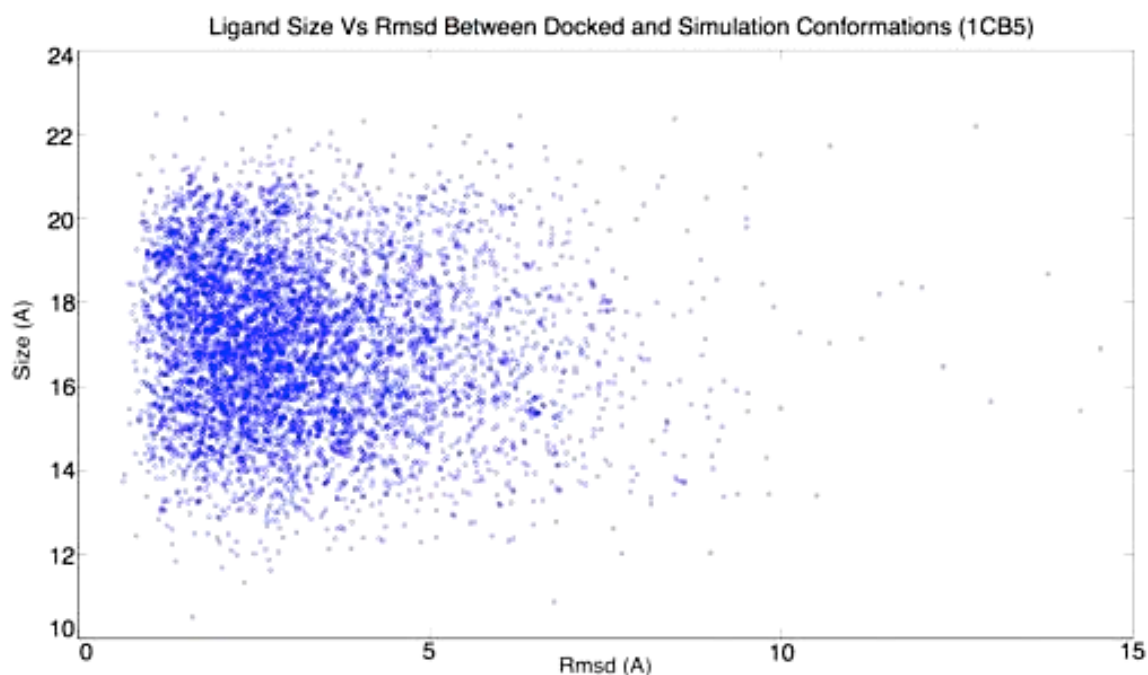
When comparing rmsd differences between the docked and simulation conformations, it can be concluded that the rmsd difference is not depended on the size of the ligand (Figure 34 and Figure 35). In the case of human T-protein, it seems true that smaller size molecules show

smaller rmsd differences between the docked and simulation binding modes (lower left quarter of Figure 34), however it is also the smaller size molecules that have the highest rmsd differences (lower right quarter of Figure 34). Interestingly, the largest molecules seem to have quite similar docked and simulation binding modes and smaller rmsd differences than smaller molecules (upper left quarter of Figure 34). This might be thanks to the well-defined tunnel-like binding site of the human T-protein. Since larger molecules wouldn't be as free as the smaller ones, they might have been restricted to fewer binding modes in both dockings and simulations.

While we can still find patterns regarding the relation between the rmsd difference and the ligand size in the case of human T-protein, a plausible relation between the rmsd difference and ligand size is not present for the candidate ligands of human bleomycin hydrolase (Figure 35). Larger rmsd values were obtained among both the smaller and the larger molecules (right half of Figure 35). In addition, the smallest rmsd values don't seem to be specific for smaller molecules, unlike in the case of human T-protein. The large solvent-exposed binding site of the human bleomycin hydrolase enables a larger area to explore for both smaller and larger molecules, thus allowing a wider variety of binding modes both in dockings and simulations. Therefore, it can be concluded that the difference between the docked and simulation binding modes is depended on the receptor-binding site rather than the ligand size.

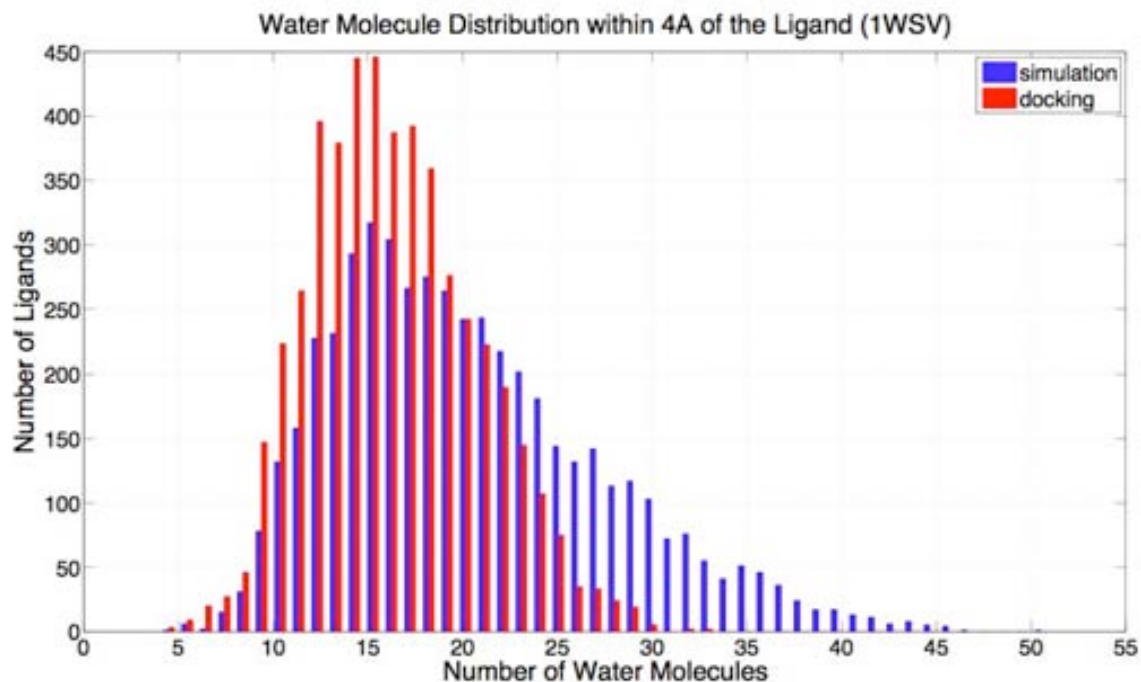


**Figure 34:** Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) versus the ligand size for human T-protein.

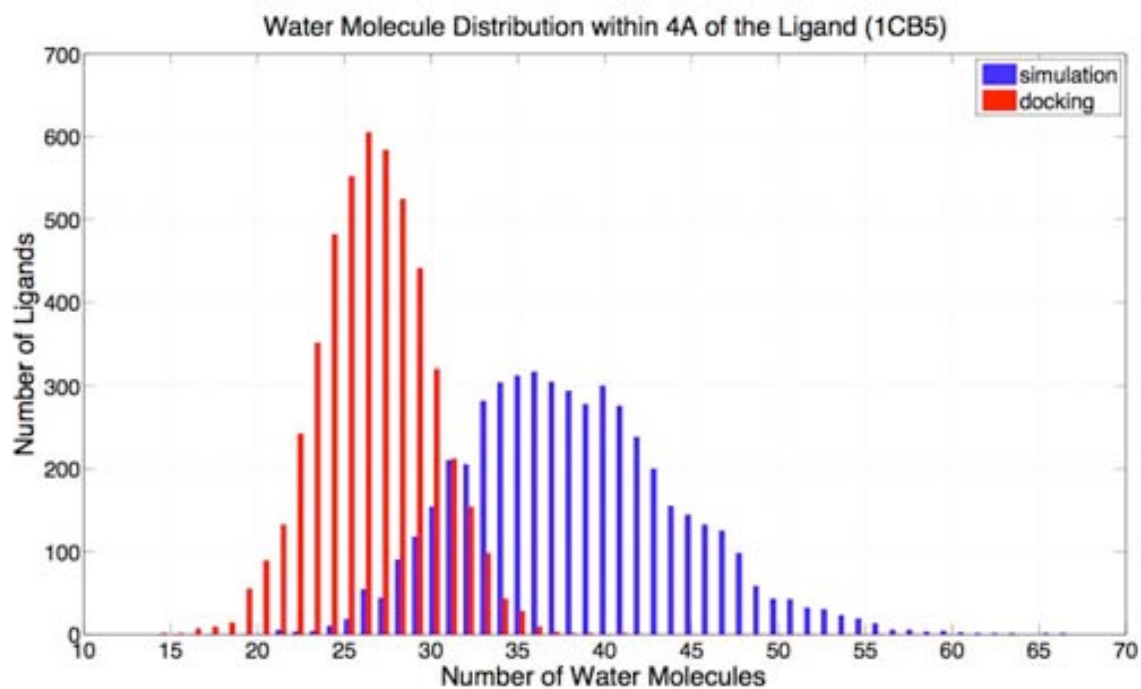


**Figure 35:** Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) versus the ligand size for human bleomycin hydrolase.

Since the docked conformations were used as the initial state inputs for the LIE simulations, the docked ligand-protein complexes were solvated, enabling an analysis of the relaxation of water molecules during the LIE simulations. The comparison of number of water molecules around the ligands at the beginning (solvated docked complexes) and end of the LIE simulations shows that the degree of solvation of the ligand increases during the simulation (Figure 36 and Figure 37). For human T-protein, since the binding site was a narrow pocket, the number of water molecules that could be in contact with the bound ligand were overall less than the number of water molecules that could surround the ligand in the large and exposed binding pocket of the human bleomycin hydrolase. In the case of both target proteins, the number of water molecules 4 Å around the bound ligand is less in general for the initial docked and solvated structure than after the LIE simulations. This is not only due to water relaxation in the binding site but also a reflect of the flexibility of both protein and ligand in the LIE simulations, accommodating (solvent-influenced) configurations that are not accessible to the docking procedure.



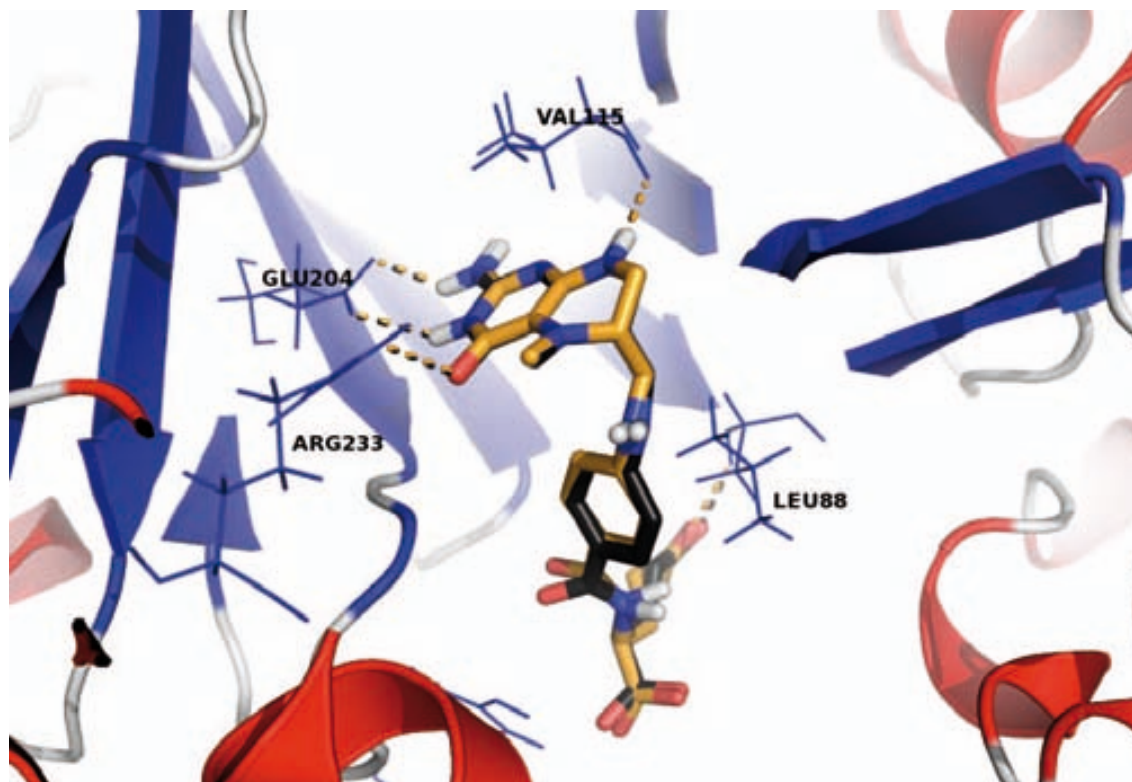
**Figure 36:** The distribution of water molecules within 4 Å of the ligands at the end of docking experiments and LIE simulations for human T-protein.



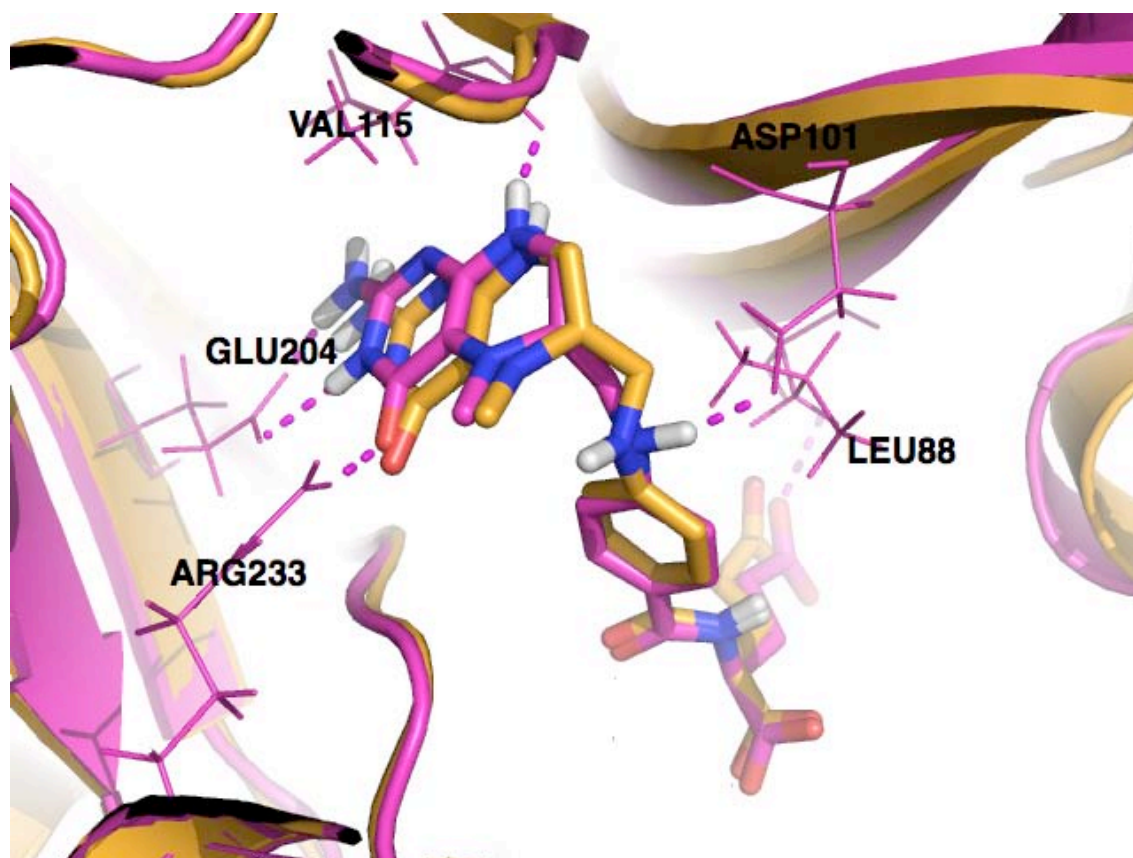
**Figure 37:** The distribution of water molecules within 4 Å of the ligands at the end of docking experiments and LIE simulations for human bleomycin hydrolase.

As a control case, we applied all steps of the workflow on the competitive inhibitor (5-CH<sub>3</sub>-H<sub>4</sub>-folate) of human T-protein and compared the results with the PDB structure. Docking 5-

CH<sub>3</sub>-H<sub>4</sub>-folate with AutoDock Vina to the binding site gave 9 poses and the one with the lowest binding free energy was chosen without visual inspection. The calculated energy for this pose by AutoDock Vina was -10.9 kcal/mol and it ranked 15<sup>th</sup> among 58699 molecules of VSL-1 when ordered according to AutoDock Vina score. Ordering according to normalized consensus score put this binding pose of 5-CH<sub>3</sub>-H<sub>4</sub>-folate to 95<sup>th</sup> place; therefore it was among the 5000 molecules selected for LIE simulations. The binding free energy calculated with LIE was -8.76 kcal/mol, however it would rank 3067 among 5000 if ordered according to LIE calculated binding free energies. The binding mode predicted by AutoDock Vina is very close to the PDB structure mode and the same hydrogen bonding pattern can be seen (Figure 38). Both poses make the double hydrogen bond with the side chain of Glu204, one hydrogen bond with Arg233 and one hydrogen bond with Leu88. The last configuration from the LIE simulations displays a very similar pose to both the docked conformation and the reference PDB structure (Figure 39 and Figure 40). The hydrogen bonding pattern for the final LIE conformation is the same as the reference and docked binding modes, with an additional hydrogen bond with Asp101. Even though the LIE binding mode is quite close to the actual binding mode, it is not as accurate as the docked one. On the basis of these results, one could be tempted to say that the LIE approach is inferior to standard docking with AutoDock Vina. However, this will be case dependent, as the comparison of methods shown previously suggests. In addition, the information on the dynamics of the system (i.e. stability in the binding site) might help deciding on a short list of candidates even when considering the docking scores as first filter rather than the LIE free energies. A similar control experiment couldn't be done with human bleomycin hydrolase since a structure with a known ligand is not available.



**Figure 38:** The competitive inhibitor 5-CH<sub>3</sub>-H<sub>4</sub>-folate in the binding site. Yellow colored binding mode is taken from the PDB structure 1WSV, while black colored binding mode is the lowest binding free energy pose generated by AutoDock Vina. Dashed lines show the corresponding hydrogen bonds.



**Figure 39:** The competitive inhibitor 5-CH<sub>3</sub>-H<sub>4</sub>-folate in the binding site. Yellow colored binding mode is taken from the PDB structure 1WSV, while magenta colored binding mode is the last configuration from the LIE simulations. Dashed magenta lines show the hydrogen bonds formed between the ligand and the target at the end of LIE simulations. Hydrogen bonds for the reference PDB structure 1WSV are not shown.



**Figure 40:** Configurations of the competitive inhibitor 5-CH<sub>3</sub>-H<sub>4</sub>-folate superimposed. Yellow molecule is taken from the PDB structure 1WSV, black molecule is the lowest binding free energy pose generated by AutoDock Vina and magenta molecule is the last configuration from the LIE simulations.

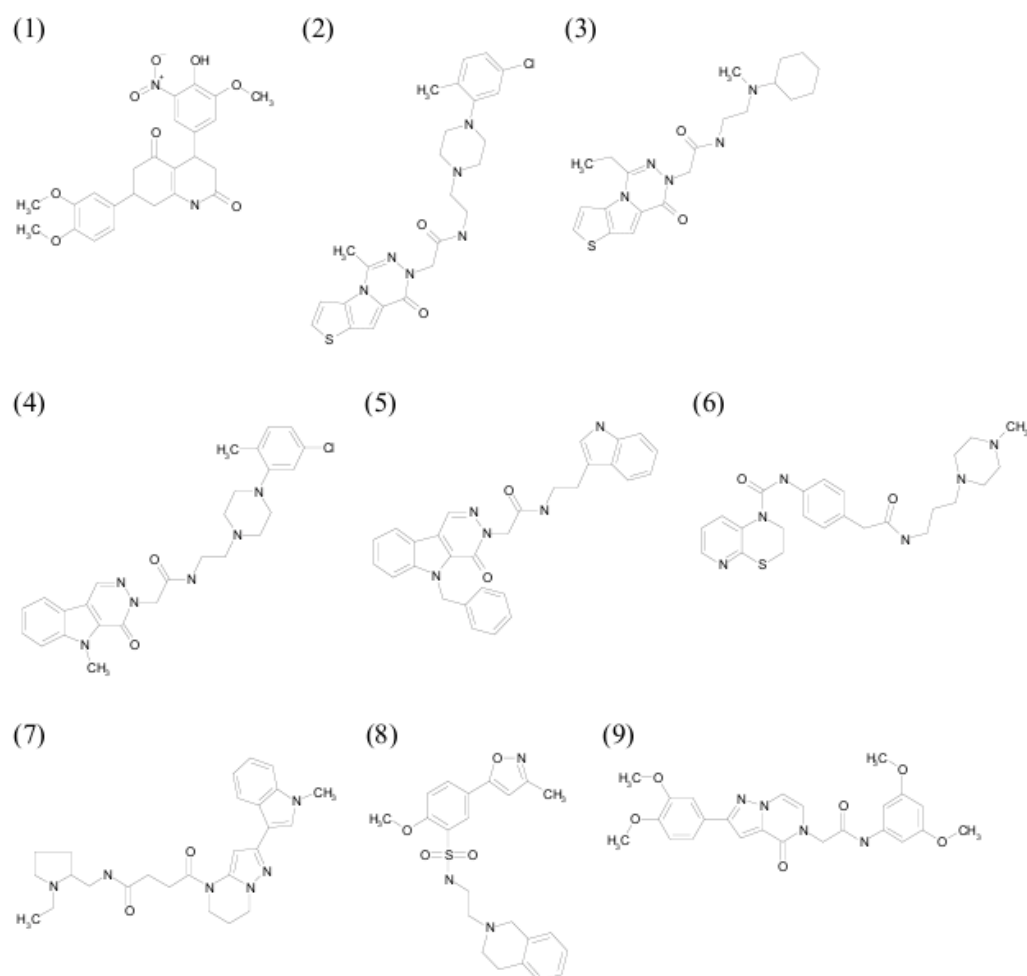
The workflow we present here is a combination of different methodologies, where information has to be passed between different applications. While developing the workflow interoperability, special care had to be taken when working with small molecule coordinate files and converting between different formats used by different programs and when re-adding protons after a step performed by a tool working with a united-atoms approach.

For human bleomycin hydrolase, we selected 40 molecules. The selection process was done in two parts. In the first part, the 4960 molecules evaluated with LIE were visually analyzed according to their docked conformations. A set of 100 compounds was selected based on their hydrogen bonding patterns with the binding site residues and the abundance of other noncovalent interactions. Then these compounds were ordered according to their binding free energy calculated with LIE. Compounds with a binding free energy smaller than -10 kcal/mol made it to the final selection set. The final selection set has 14 compounds chosen with this criterion. The second part of the selection process was done taking only into account the binding free energies calculated with LIE. The 4960 molecules were ordered based on their binding free energy and compounds with a value smaller than -15 kcal/mol were included in the final selection set. This second criterion contributed 26 compounds to the final selection. From this set of 40 molecules, 9 were available from a single vendor and they were purchased for experimental testing. Figure 41 shows the 2D structures of the experimentally tested molecules, and the selection criteria and binding free energies calculated with the LIE method are summarized in **¡Error! No se encuentra el origen de la referencia..**

Compound no.	Selection criteria	LIE Energy (kcal/mol)	Rmsd between docking and LIE (Å)
1	LIE	-16	5.2
2	LIE	-17	5.4
3	Visual	-11	2.3
4	LIE	-16	1,6
5	LIE	-16	3,5
6	LIE	-16	5.7
7	LIE	-16	2.3
8	LIE	-15	2.8
9	LIE	-16	4.4

**Table 13:** The list of selection criteria fulfilled by the chosen molecules for testing against human bleomycin hydrolase.





**Figure 41:** The molecules selected for experimental testing of inhibition of human bleomycin hydrolase.

The experimental testing of the chosen molecules was carried out by Dr. Lionel Costenaro. The fluorogenic substrate L-citrulline 7-amido-4-methylcoumarin hydrobromide (H-Cit-AMC·HBr; Bachem) was used to assay the protease activity of recombinant hBH *in vitro* and its inhibition by the compounds. Cleavage reaction mixture contained 50 nM hBH, substrate (dissolved in water) and potential inhibitor (dissolved in 100% DMSO), in 140 mM KCl, 100 mM Tris·HCl pH 7.5, in a total volume of 100  $\mu$ l. Seven substrate concentrations (*S*) between 10 and 1500  $\mu$ M and four inhibitor concentrations (*I*) between 1 and 100  $\mu$ M were typically used to assay the effect of one potential inhibitor. The final DMSO concentration was 2%. After 30 min. of incubation with the potential inhibitor, the reactions were initiated by the addition of the substrate and performed at 30°C in 96-well microlitre plates (white wells, black frame). All conditions were done in triplicates. The fluorescence of the liberated product (AMC) was monitored for 6h (excitation at a wavelength of 355 nm and emission at 460 nm) by a plate reader Victor III (Perkin Elmer). Fluorescence intensities were transformed to AMC amount (pmol) using a standard curve of 7-amido-4-methylcoumarin (Bachem).

In inhibition experiments, the cleavage reaction mixture contained 2% of DMSO, regardless the potential inhibitor concentrations. To assess the effect of DMSO on hBH kinetics, we performed a control experiment with 2% DMSO, but without inhibitor. The differences

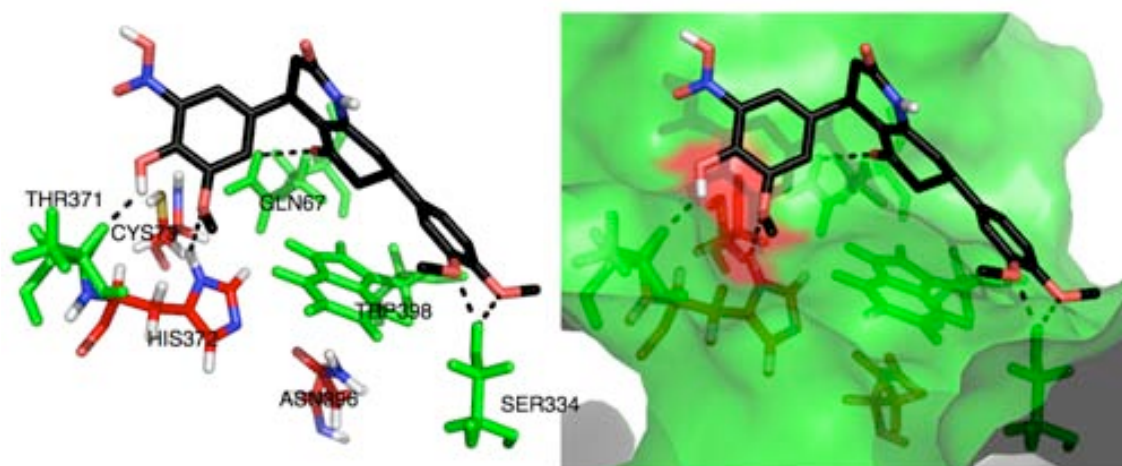
between  $V/V_{max}$  values with and without 2% DMSO were similar to the variation between experimental replicates for all substrate concentrations. This clearly shows that 2% DMSO does not influence hBH activity (Table 14).

Out of nine compounds, Compounds **3**, **6** and **8** did not inhibit the hydrolysis of H-Cit-AMC by hBH and Compound **4** was not soluble at a concentration sufficient to be tested. Values of  $K_m$  in presence of inhibitor were similar to that of hBH alone, except for compound **7**, which had the highest determined  $K_i$ . The compounds showed different modes of inhibition: competitive, noncompetitive or mixed (Table 14). The derived inhibition constants  $K_i$  ranged from 30 to 356  $\mu\text{M}$ .

No	Mode of inhibition	$K_m$ [ $\mu\text{M}$ ]	$K_i$ [ $\mu\text{M}$ ]
Control	n.a.	182 $\pm 12$	n.a.
1	Noncompetitive	190 $\pm 10$	30 $\pm 2$
2	Mixed	195 $\pm 14$	89 $\pm 11$
3	No inhibition	163 $\pm 9$	n.a.
4	n.a.	n.a.	n.a.
5	Noncompetitive	212 $\pm 16$	263 $\pm 24$
6	No inhibition	176 $\pm 16$	n.a.
7	Competitive	138 $\pm 7$	356 $\pm 155$
8	No inhibition	199 $\pm 10$	n.a.
9	Mixed	221 $\pm 14$	127 $\pm 48$

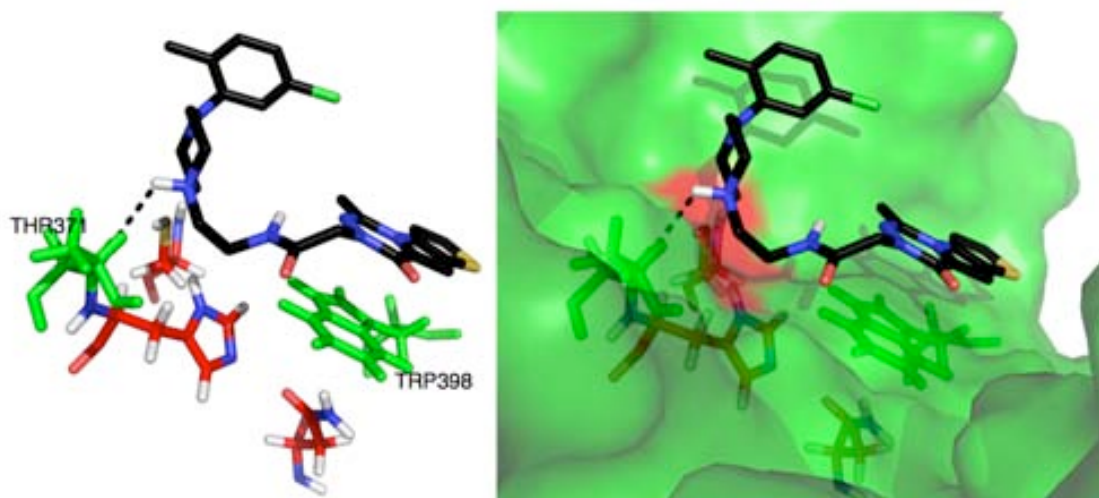
**Table 14:** Observed activities of the candidate compounds for human bleomycin hydrolase.

When the position predicted by docking of Compound **1** in the binding site is examined, it is seen that the compound is stretched along the binding site contacting the catalytic residues and making hydrogen bonds with residues Gln67, Ser334, Thr371 and His372 (Figure 42). There is also a  $\pi$ - $\pi$  interaction between Compound **1** and the indole ring of Trp398.



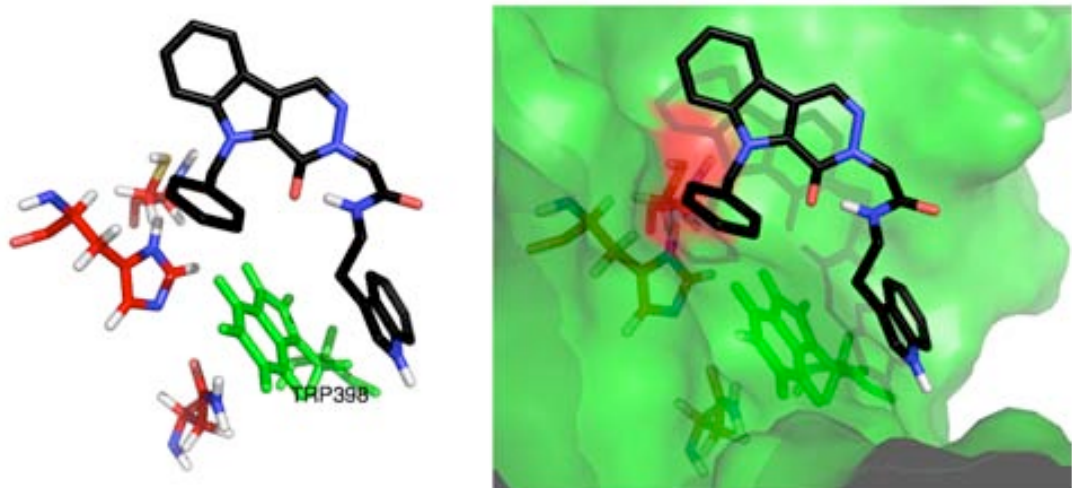
**Figure 42:** Binding mode of Compound **1** as predicted by docking. Compound **1** is shown in black, catalytic residues Cys73, His372 and Asn396 are shown in red while other residues contacting the compound are shown in green.

Figure 43 shows the conformation of Compound **2** in the binding site predicted with docking. Compound **2** is stabilized with a hydrogen bond with Thr371 and  $\pi$ - $\pi$  interactions with Trp398.



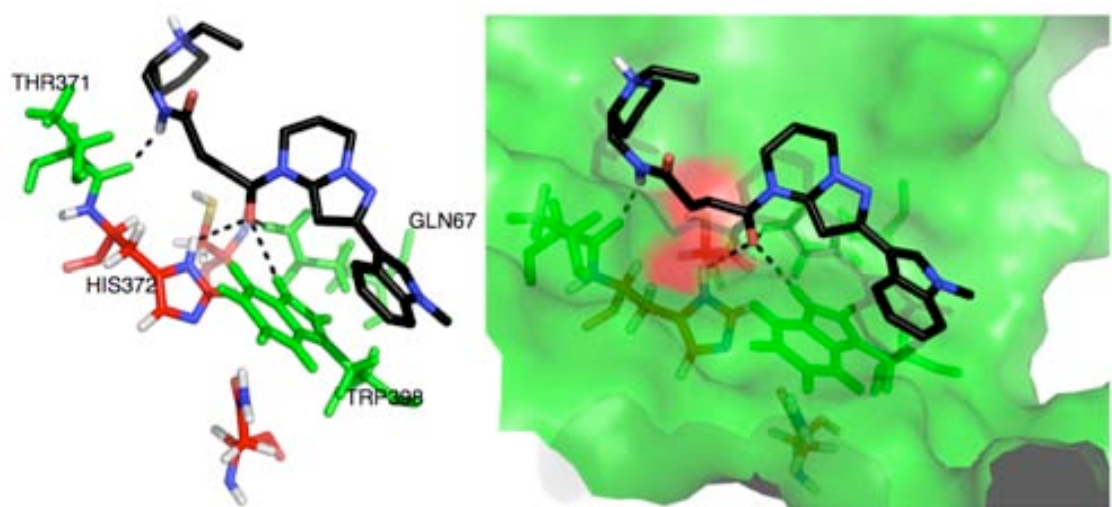
**Figure 43:** Binding mode of Compound **2** as predicted by docking.

The only interaction that Compound **5** seems to have in its docked conformation is with Trp398 (Figure 44). The indole ring of Trp398 stabilizes the ringed tail of Compound **5** through  $\pi$ - $\pi$  interactions.



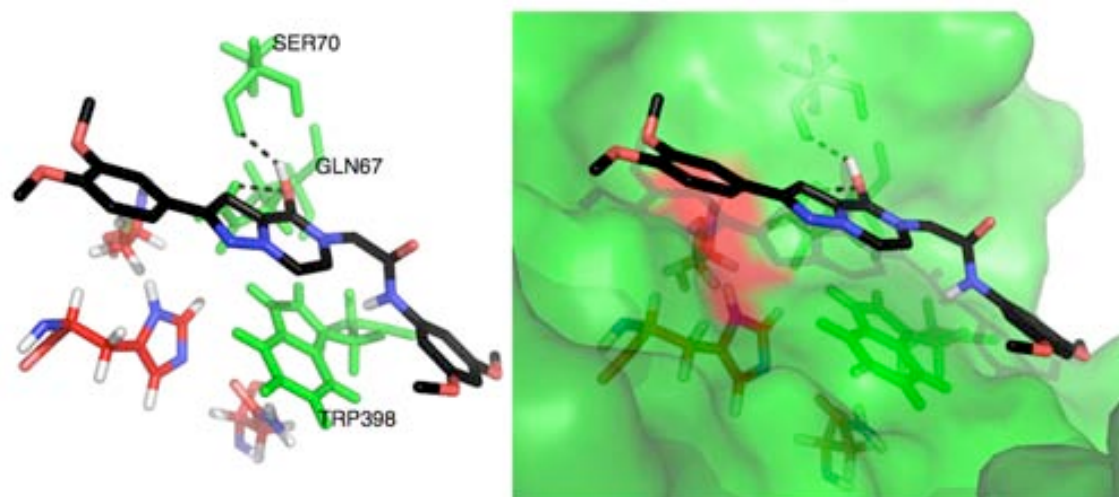
**Figure 44:** Binding mode of Compound 5 as predicted by docking.

Compound 7 is also extended along the surface-exposed binding site making hydrogen bonds with Gln67, Thr371, His372 and Trp 398 in the docked conformation (Figure 45). Trp398 also interacts with Compound 7 through  $\pi$ - $\pi$  interactions.



**Figure 45:** Binding mode of Compound 7 as predicted by docking.

In the docked conformation, Compound 9 is also extended along the binding site, making contacts with Gln67 and Ser70 (Figure 46). Like the other compounds, Compound 9 also makes  $\pi$ - $\pi$  interactions with the indole ring of Trp398.



**Figure 46:** Binding mode of Compound **9** as predicted by docking.

#### 4.3.2. Results For Human Gcase

The 3D flexible pharmacophore search with **pharma1** reduced the library size to 136252 and **pharma2** reduced it to 206428 from 2157575. Both filters managed to decrease the number of molecules to be docked by a factor of around 10. Even though no primary restrictions about the sizes of ligands were designed in either of the pharmacophore filters, large molecules (molecules with more than 20 Å for the longest distance between two atoms) were mostly filtered out in both cases. Out of the 19678 large molecules in the small molecule library, only 769 passed the filter **pharma1** and 2518 passed **pharma2**. The sizes of the ligands to be docked are important because it is known that, regardless of the scoring function used, larger molecules tend to produce better scores than smaller molecules simply because of the abundance of hypothetical interactions in the binding sites.<sup>37,85</sup>

After the docking experiments **dock1**, **dock2** and **dock3**, the poses with the lowest binding free energies calculated by the corresponding scoring function (AutoDock Vina scoring function for **dock1** and **dock2**, Surflex-Dock scoring function for **dock3**) were selected for rescoring with the four scoring functions implemented in CSCORE and subsequent consensus scoring. The scores given by each scoring function were normalized to values between 0 and 1.

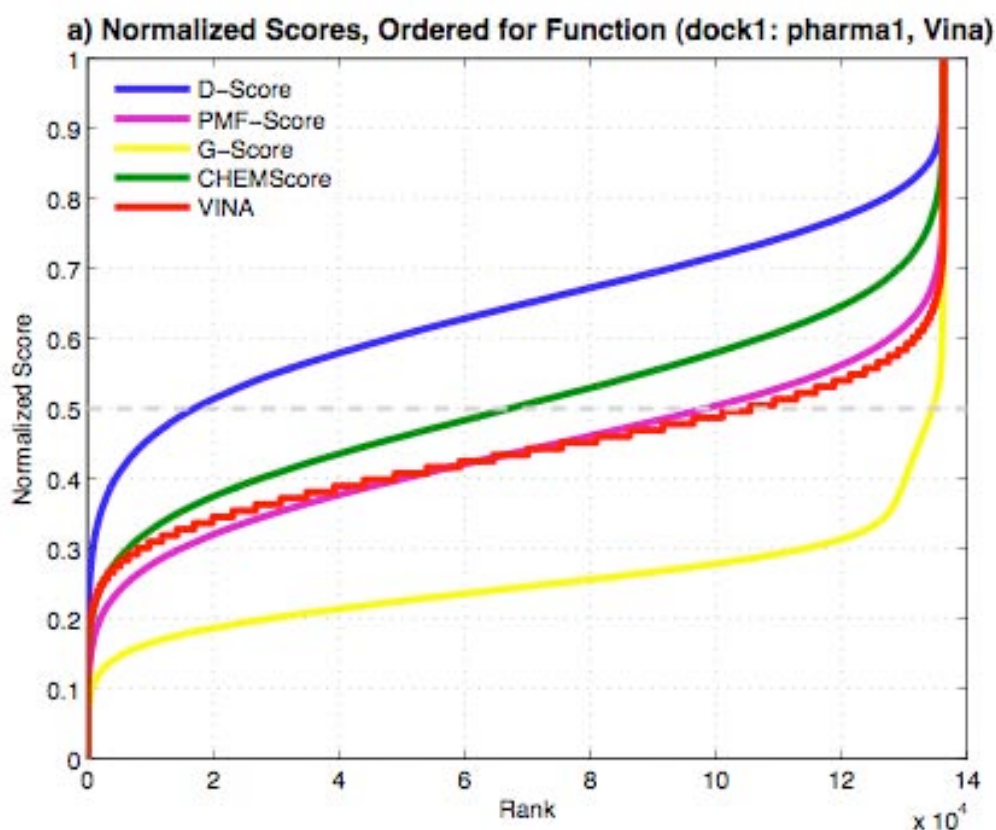
Table 15 shows Pearson's correlations coefficients between different scoring functions for the docking experiment **dock1** for the **pharma1** filtered library.

NCS	D-Score	PMFScore	G-Score	ChemScore	VinaScore
-----	---------	----------	---------	-----------	-----------

NCS	1	0.76	0.62	0.65	0.76	0.75
D-Score	0.76	1	0.48	0.42	0.49	0.33
PMFScore	0.62	0.48	1	0.17	0.13	0.36
G-Score	0.65	0.42	0.17	1	0.43	0.35
ChemScore	0.76	0.49	0.13	0.43	1	0.63
VinaScore	0.75	0.33	0.36	0.35	0.63	1

**Table 15:** Pearson's Correlation coefficients between different scoring functions and  $NCS_{99.5}$  for docking experiment **dock1**

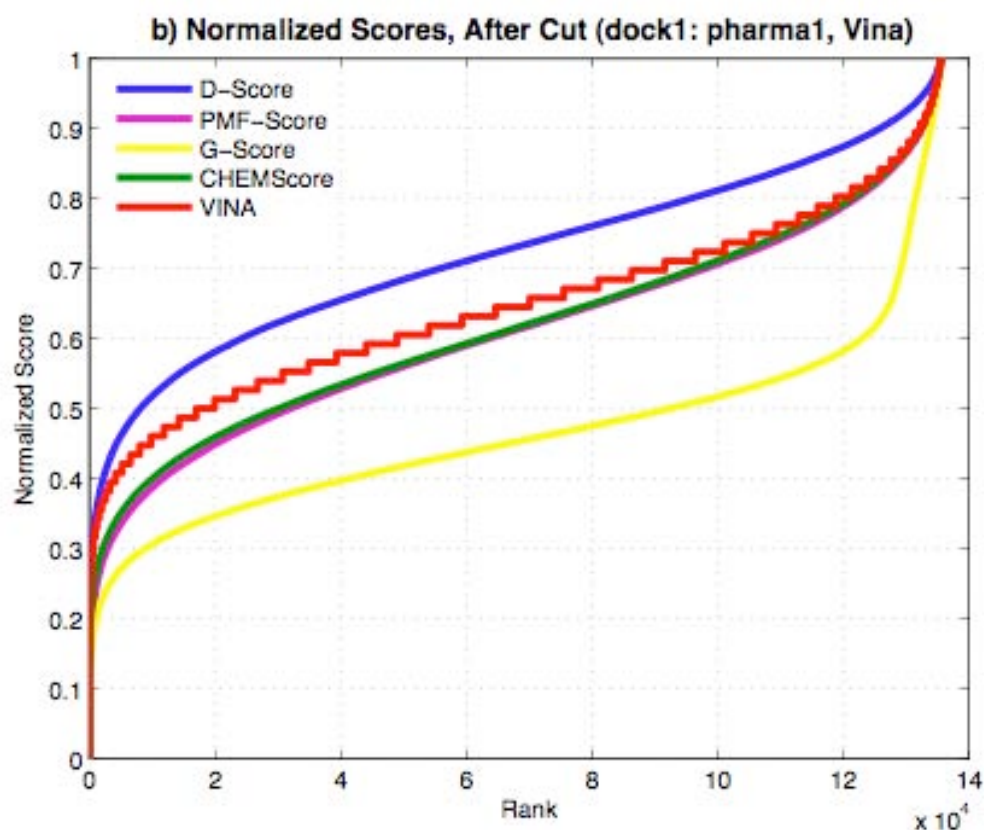
For docking experiment **dock1**, normalization of the scores output by different scoring functions to values between 0 and 1 with a truncation *cut-off* of 100% generated similarly shaped sigmoidal curves for all scoring functions (Figure 47). However, especially in the case of G-Score, the individual scoring functions showed trails of high values to largely different extents. For the scoring functions except for D-Score and ChemScore, the values at the poor scoring ends are quite different than the rest, making the distinction between the fairly well docked ligands and the fairly poor docked ones quite difficult. These values at the poorly scoring end are also problematic when combining the different scores to reach a consensus. Using a rank-by-vote strategy for consensus scoring based on scores being within the top  $n\%$  of the obtained score range, with a frequently-used "vote cut-off" of 0.5 (dashed grey line in Figure 47), as in CSCORE, results in G-Score voting for more than 99% of the molecules, and PMFScore and Vina score for around 75% of the molecules, rendering these scoring functions nearly redundant.



**Figure 47:** Normalized scores calculated without truncation with the five scoring functions Vina, PMFScore, G-Score, D-Score and ChemScore for **dock1** experiment of GCase.

Using a smaller “vote cut-off” value (around 0.3) would still not be enough because this time, all scoring functions except G-Score would be very specific and decisive and vote for a very small amount of molecules, while G-Score would still give a passing vote to more than 90% of the molecules. Summing the normalized scores,  $NCS_{100}$ , would be problematic as well because of the discrepancies between the decisiveness of individual scoring functions. For example, a compound ranked around 10,000 by D-Score receives a score that is twice that of another molecule ranked around 1,000 by D-Score, making D-Score highly decisive and sensitive against fairly poor and fairly well docked molecules. However, G-Score scores almost 100,000 of the molecules with almost the same value, showing no sensitivity except for very poorly docked molecules. Therefore, a solution to close the gap between the decisiveness and sensitivity of different scoring functions was to exclude the very poorly scoring end of each scoring function and to calculate the  $NCS$  score after the truncation.

Normalization with a truncation *cut-off* of 99.5% decreased the slope of G-Score at the poorly scoring end, increasing the overall sensitivity (Figure 48). The curves of Vina score, ChemScore and PMFScore became more similar to the curve of D-Score. Even though G-Score’s sensitivity is increased, it is still not close to the rest. To make G-Score converge with the rest would require the truncation of at least 5% of the molecules instead of 0.5%. However we decided not to diminish the number of molecules and thus it was decided to stop at 0.5% truncation and calculate the  $NCS_{99.5}$  values for each compound.



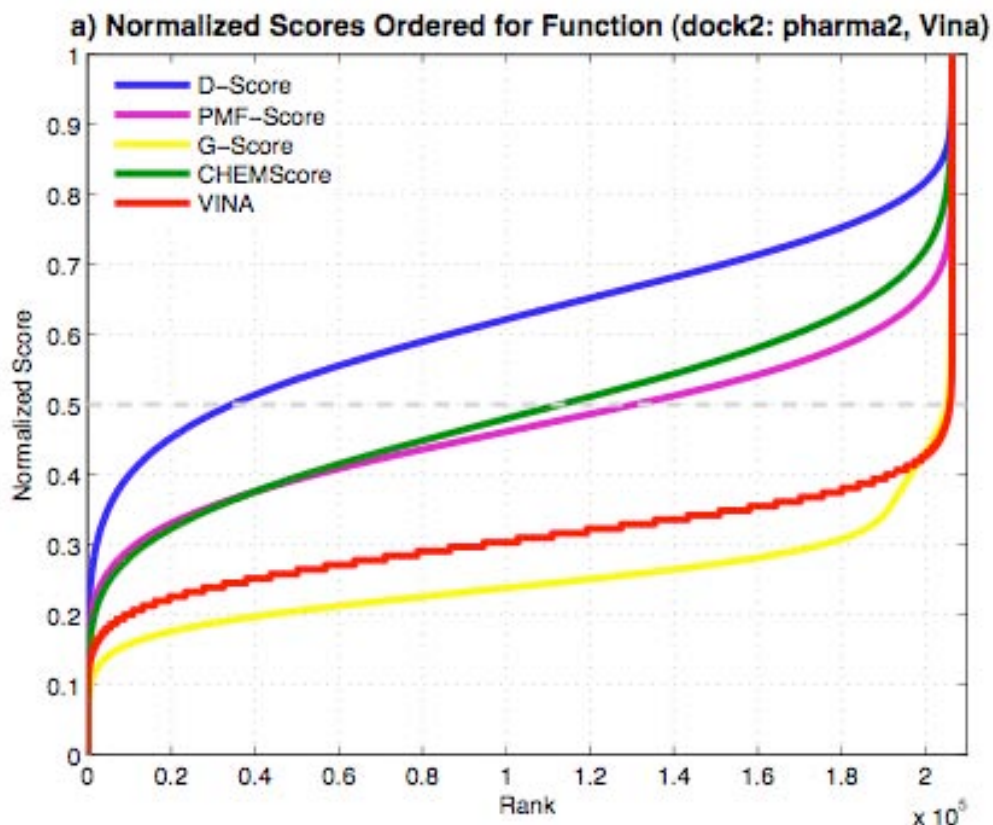
**Figure 48:** Normalized scores calculated with the second normalization procedure against compound rank (after the worst-scoring 0.5% were excluded) for **dock1** experiment of GCASE.

Even though the molecules docked in the experiment **dock2** were filtered by a different pharmacophore filter, the curves of scoring functions are not very different from those of **dock1**, with the curve for Vina score being the exception. Normalization with 100% *cut-off* (Figure 49) shows that Vina score and G-Score gave very high penalties to a few very poorly docked molecules, making the area between very well and very poor scoring ends quite flattened. Ranking the molecules according to their  $NCS_{100}$  would be problematic for the experiment **dock2** as well, since D-Score is sensitive to poorly docked molecules while almost 99% of the molecules would be classified as well-docked and only a small amount of molecules would receive a high normalized score by G-Score and Vina score.

	NCS	D-Score	PMFScore	G-Score	ChemScore	VinaScore
NCS	1	0.77	0.64	0.67	0.77	0.75
D-Score	0.77	1	0.51	0.43	0.49	0.33
PMFScore	0.64	0.51	1	0.23	0.19	0.38
G-Score	0.67	0.43	0.23	1	0.42	0.35
ChemScore	0.77	0.49	0.19	0.42	1	0.65
VinaScore	0.75	0.33	0.38	0.35	0.65	1

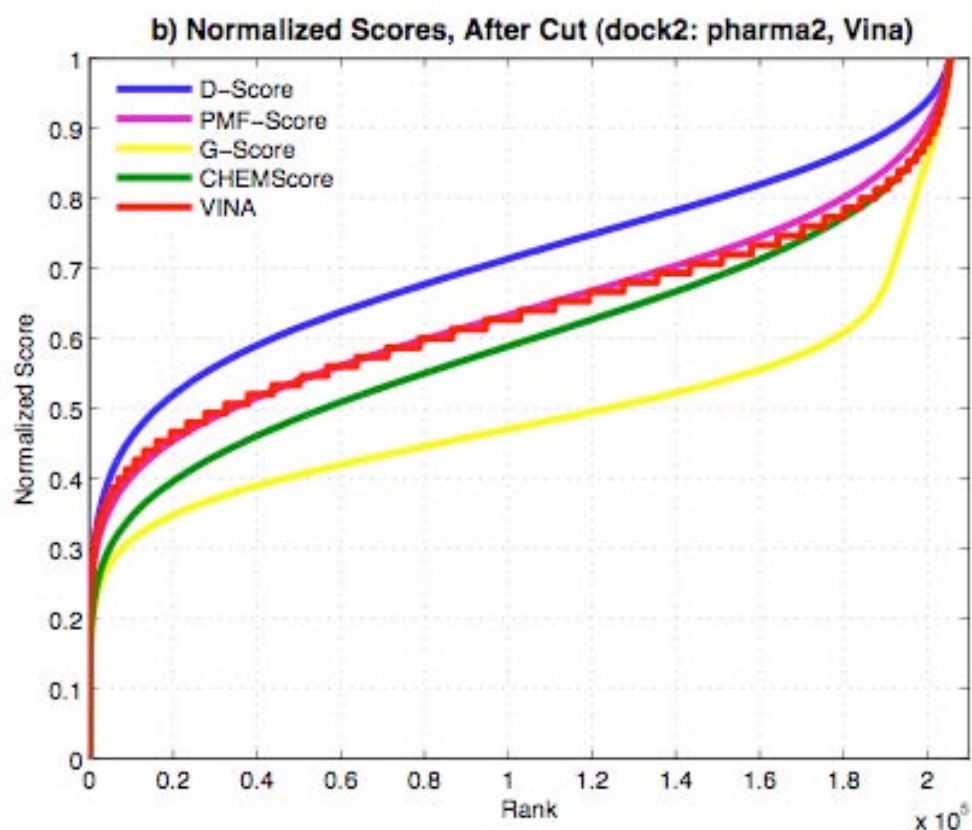
**Table 16:** Pearson's Correlation coefficients between different scoring functions and  $NCS_{99.5}$  for docking experiment **dock2**.





**Figure 49:** Normalized scores calculated without truncation with the five scoring functions Vina, PMFScore, G-Score, D-Score and ChemScore for **dock2** experiment of GCASE.

Using a truncation *cut-off* of 99.5% for normalization brought the curves of different scoring functions close to each other as in the case of **dock1** (Figure 50). While cutting the poor-scoring end brought VINA score on the same order as D-Score, ChemScore and PMFScore, G-Score's improvement was not as remarkable. G-Score would need more molecules to be excluded to converge, and thus would reduce the number of molecules even further. However, even with a truncation cut-off of 99.5%, it would still vote for 50% of the molecules with a "vote cut-off" of 0.5, therefore it was decided to stop truncation at 0.5% and calculate the  $NCS_{99.5}$  values to proceed to rank the molecules. Even though the number of molecules docked in experiment **dock2** was twice as large, all correlation values of **dock2** are very similar to those of **dock1** (Table 16).



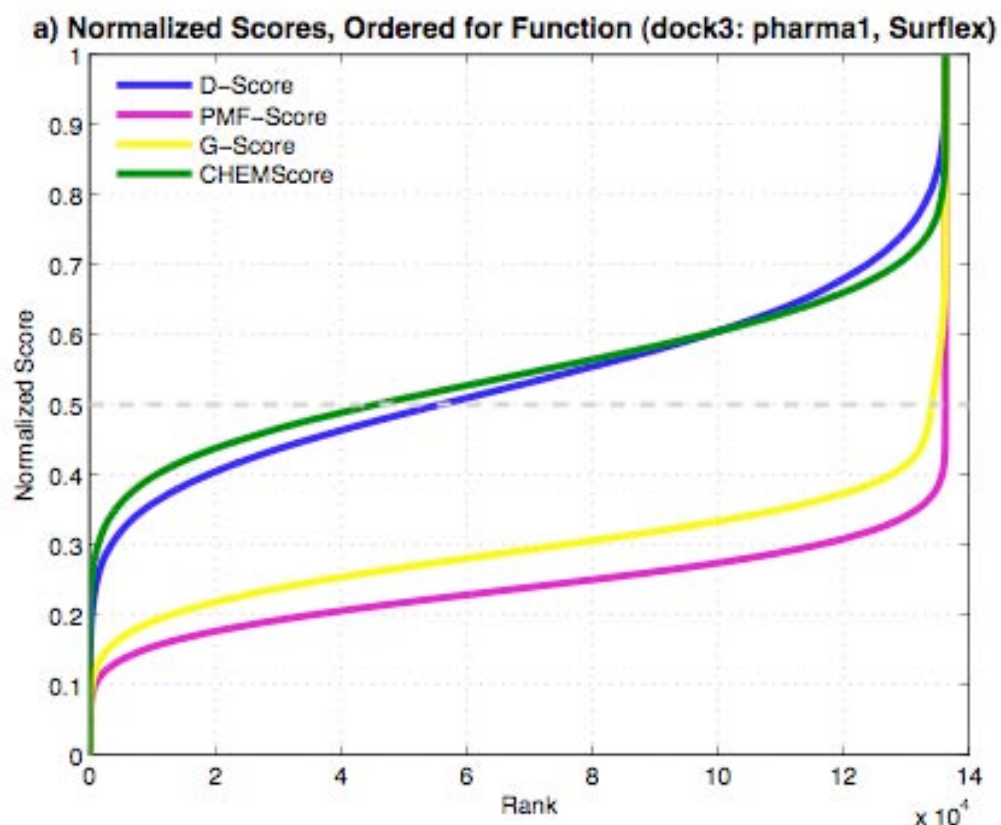
**Figure 50:** Normalized scores calculated with the second normalization procedure against compound rank (after the worst-scoring 0.5% were excluded) for **dock2** experiment of GCASE.

Correlation values of **dock3** are overall better than those of **dock1** and **dock2** (Table 17). PMFScore is better correlated with the other scoring functions, while better correlation values between  $NCS_{99.5}$  and the scoring functions is also observed. Since **dock1** and **dock3** experiments were done with the same set of molecules, the differences between the correlation values enables us to make comparisons between **dock1** done with AutoDock Vina and **dock3** done with Surflex-Dock. In **dock3**, individual scoring functions, especially PMFScore, show improved correlations with each other and also with  $NCS_{99.5}$ . This shows that Surflex-Dock was able to find binding modes that were concurred more consistently and coherently by the additional scoring functions.

	NCS	D-Score	PMFScore	G-Score	ChemScore
NCS	1	0.9	0.71	0.76	0.79
D-Score	0.9	1	0.46	0.73	0.66
PMFScore	0.71	0.46	1	0.28	0.45
G-Score	0.76	0.73	0.28	1	0.44
ChemScore	0.79	0.66	0.45	0.44	1

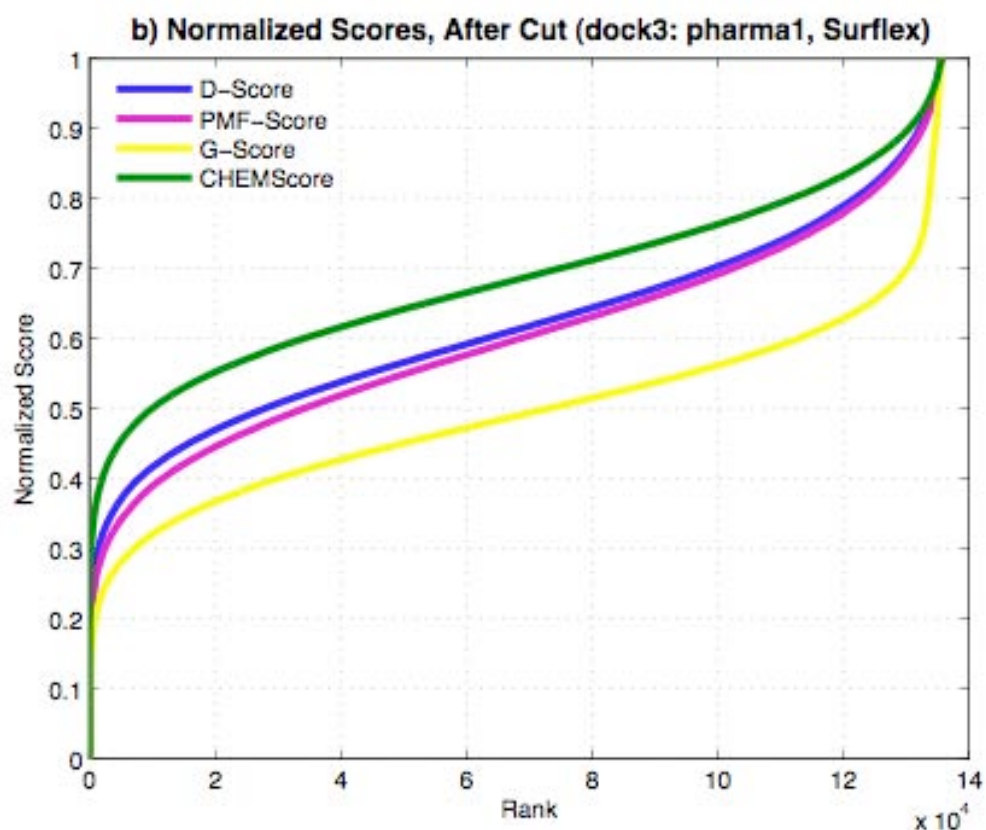
**Table 17:** Pearson's Correlation coefficients between different scoring functions and  $NCS_{99.5}$  for docking experiment **dock3**.

When normalized without any truncation, the four scoring functions show two patterns for the distributions of normalized values for the docking experiment **dock3** (Figure 51). While D-Score and ChemScore were equally decisive and would vote for 45-50% of the molecules with a vote cut-off of 0.5 (dashed grey line in Figure 51), both G-Score and PMF-Score scored few molecules with very high penalties, decreasing the overall sensitivity. These scoring functions would vote in favor of almost all of the molecules with 0.5 vote cut-off, making no contribution to the overall ranking.



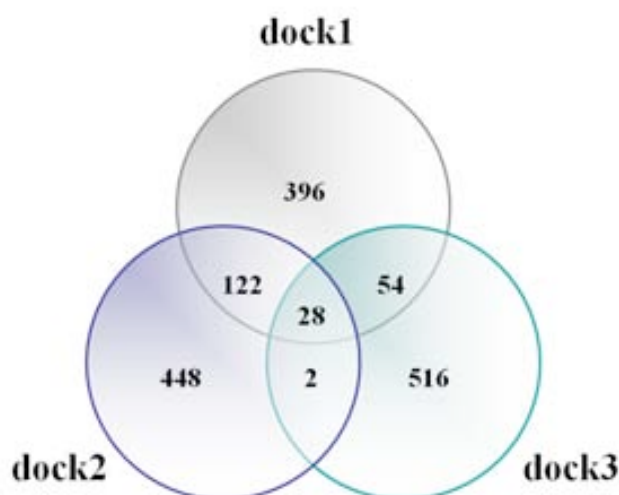
**Figure 51:** Normalized scores calculated without truncation with the four scoring functions PMFScore, G-Score, D-Score and ChemScore for **dock3** experiment of GCase.

On the other hand, with a small truncation of 0.5% of the molecules, the curves of all four scoring functions could be brought to similar sensitivity levels. G-Score again, like in the case of **dock1** and **dock2**, couldn't converge as good as the rest of the curves (Figure 52).



**Figure 52:** Normalized scores calculated with the second normalization procedure against compounds rank (after the worst-scoring 0.5% were excluded) for **dock3** experiment of GCase.

After the calculation of  $NCS_{99.5}$  values for each compound from each docking experiment, the compounds were ranked according to their  $NCS_{99.5}$ . From the three docking experiments, the compounds ranking in the top 600 were selected for the following step of binding free energy estimation with LIE simulations, adding up to 1800 compounds in total. There were some intersecting molecules, i.e. 178 molecules were in the top 600 of any two docking experiments and 28 were in the top 600 of all three docking experiments. However, repeating molecules were not reduced to only one conformation, all docked conformations of the repeating molecules from different docking experiments were included in the simulations. Each conformation was kept as input for the molecular dynamics simulations because each conformation corresponds to a different starting point and different starting points may affect the outcome of short molecular dynamics simulations dramatically. Therefore, the number of unique molecules was 1566; 1360 molecules with a single conformation, 178 molecules with two conformations and 28 molecules with three conformations, adding up to 1800 conformations in total (Figure 53).



**Figure 53:** Numbers of intersecting molecules chosen by more than one docking experiment in the top 600.

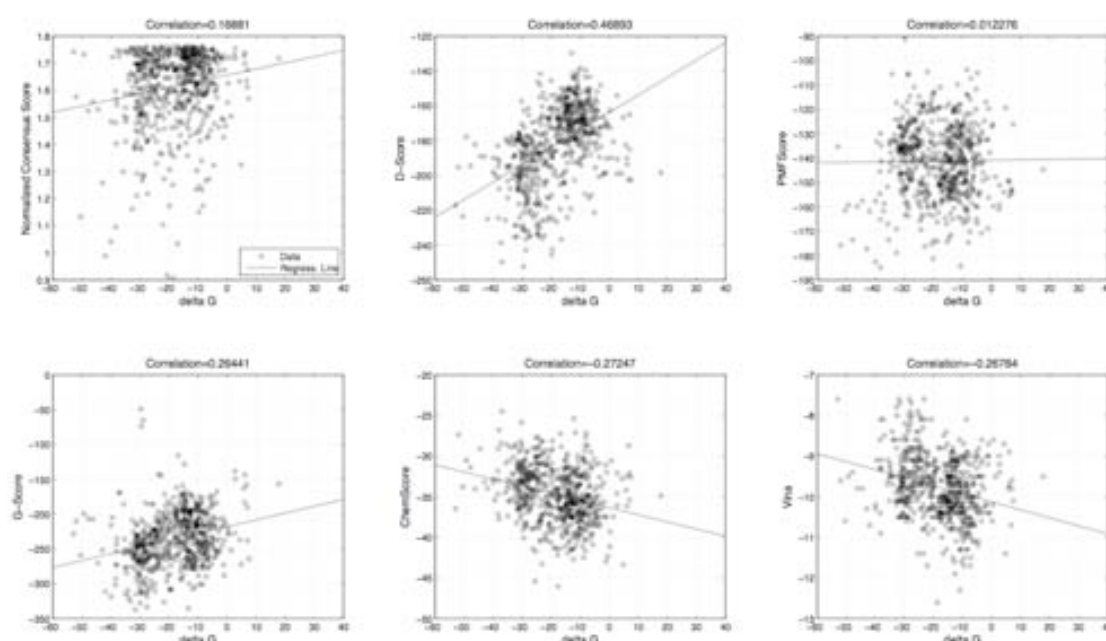
For the LIE simulations, the same force-field parameters and ligand topologies as in the case of human T-protein and human bleomycin hydrolase were used. The center and the size of the solvation sphere were again decided as the former examples. The same solvation sphere was used for the simulations of all selected ligands from the three docking experiments. Since the defining factor for the size of the solvation sphere is the maximum size of a compound and the largest molecule of the final set in its docked conformations was 24 Å, the radius of the solvation sphere was set to 27 Å.

To have a neutral environment for both the free and bound states, the following titratable residues were left charged: Arg120, Asp127, Arg131, Glu152, Asp153, Lys157, Lys186, Lys194, Glu235, Asp282, Asp283, Arg285, His311, Asp315, Glu340, Lys346, Glu349, Arg353, Arg359, Glu388, Arg395 and Asp399. The remaining charged residues were neutralized because they were out of the solvation sphere. 1566 molecules with 1800 different conformations were simulated for binding free energy calculation with LIE method.

The 600 selected molecules from the docking experiment **dock1** have been evaluated for the correlations between the values assigned by individual scoring functions and the binding free energies calculated with the LIE method (Table 18). As expected (from the analysis with the larger set of molecules given in Table 15), all individual scoring functions gave low correlation values to the normalized consensus score, *NCS*. However, the correlations among individual scoring functions were not significant for the selected 600 molecules from **dock1** experiment. The only scoring function that showed significant correlation with LIE energies for experiment **dock1** was D-Score (Figure 54).

	NCS	D-Score	PMF Score	G-Score	Chem Score	Vina	LIE
NCS	1	0.33	0.36	0.23	0.24	0.26	0.17
D-Score	0.33	1	0.08	0.12	-0.27	-0.55	0.47
PMFScore	0.36	0.08	1	-0.27	-0.39	-0.03	0.01
G-Score	0.23	0.12	-0.27	1	-0.1	-0.28	0.26
ChemScore	0.24	-0.27	-0.39	-0.1	1	0.26	-0.27
Vina	0.26	-0.55	-0.03	-0.28	0.26	1	-0.27
LIE	0.17	0.47	0.01	0.26	-0.27	-0.27	1

**Table 18:** Pearson's Correlations of scoring functions with each other and with LIE energies for the selected molecules from **dock1** experiment.

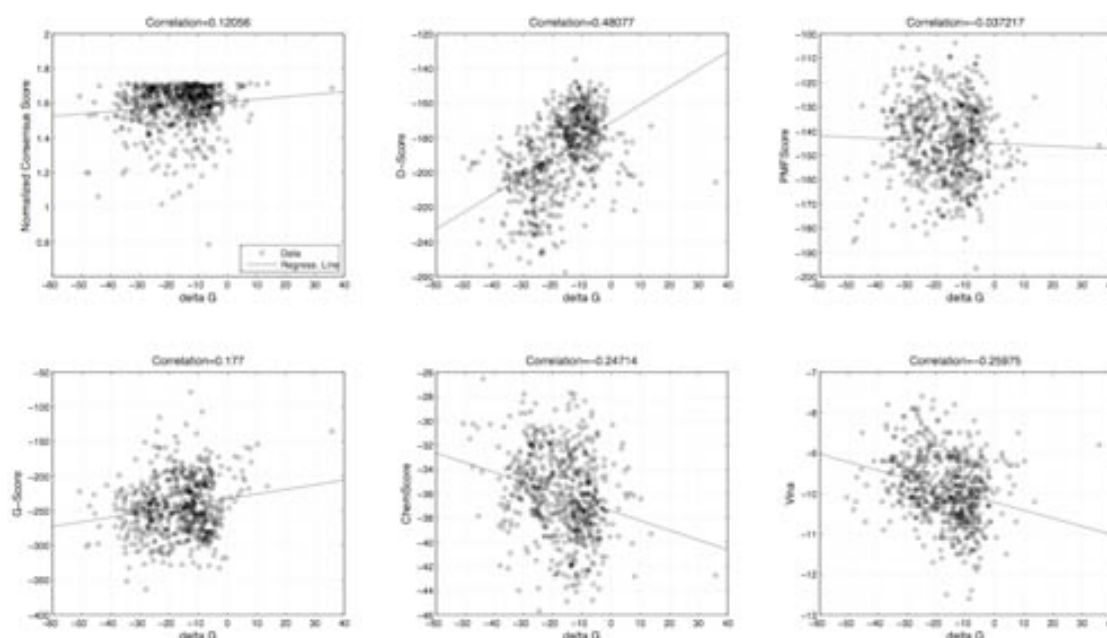


**Figure 54:** Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for **dock1**.

In the case of experiment **dock2**, which employed a different pharmacophore filter from **dock1**, the results were not very different (Table 19). Again D-Score was the only scoring function with a significant correlation of 0.48 with LIE values (Figure 55).

	NCS	D-Score	PMF Score	G-Score	Chem Score	Vina	LIE
NCS	1	0.27	0.38	0.2	0.15	0.22	0.12
D-Score	0.27	1	0.07	0.06	-0.33	-0.56	0.48
PMFScore	0.38	0.07	1	-0.3	-0.46	0.05	-0.04
G-Score	0.2	0.06	-0.3	1	0.03	-0.35	0.18
ChemScore	0.15	-0.33	-0.46	0.03	1	0.09	-0.25
Vina	0.22	-0.56	0.05	-0.35	0.09	1	-0.26
LIE	0.12	0.48	-0.04	0.18	-0.25	-0.26	1

**Table 19:** Pearson's Correlations of scoring functions with each other and with LIE energies for the selected molecules from **dock2** experiment.



**Figure 55:** Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for **dock2**.

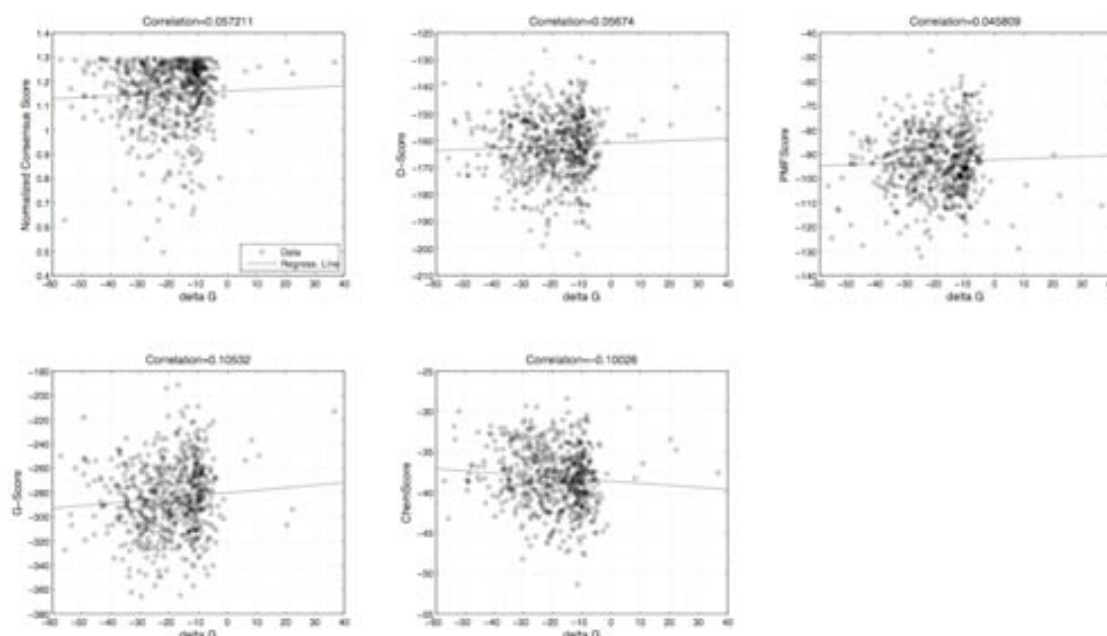
However, for the docking experiment **dock3**, a different story can be told (Table 20). First, individual scoring functions show better overall correlation with the normalized consensus scores, *NCS*. Second, except the correlation between D-Score and G-Score, individual scoring functions don't seem to be associated with each other. Especially PMFScore is significantly negatively-correlated with all remaining scoring functions. And lastly, none of the individual scoring functions are correlated with LIE energies (Figure 56).

	NCS	D-Score	PMF Score	G-Score	Chem Score	LIE
NCS	1	0.53	0.23	0.38	0.15	0.06
D-Score	0.53	1	-0.35	0.4	-0.33	0.06
PMFScore	0.23	-0.35	1	-0.43	-0.46	0.05
G-Score	0.38	0.4	-0.43	1	0.03	0.11
ChemScore	0.51	0.07	-0.22	-0.04	1	-0.1
LIE	0.06	0.06	0.05	0.11	-0.25	1

**Table**

Pearson's Correlations of scoring functions with each other and with LIE energies for the selected molecules from **dock3** experiment.

**20:**



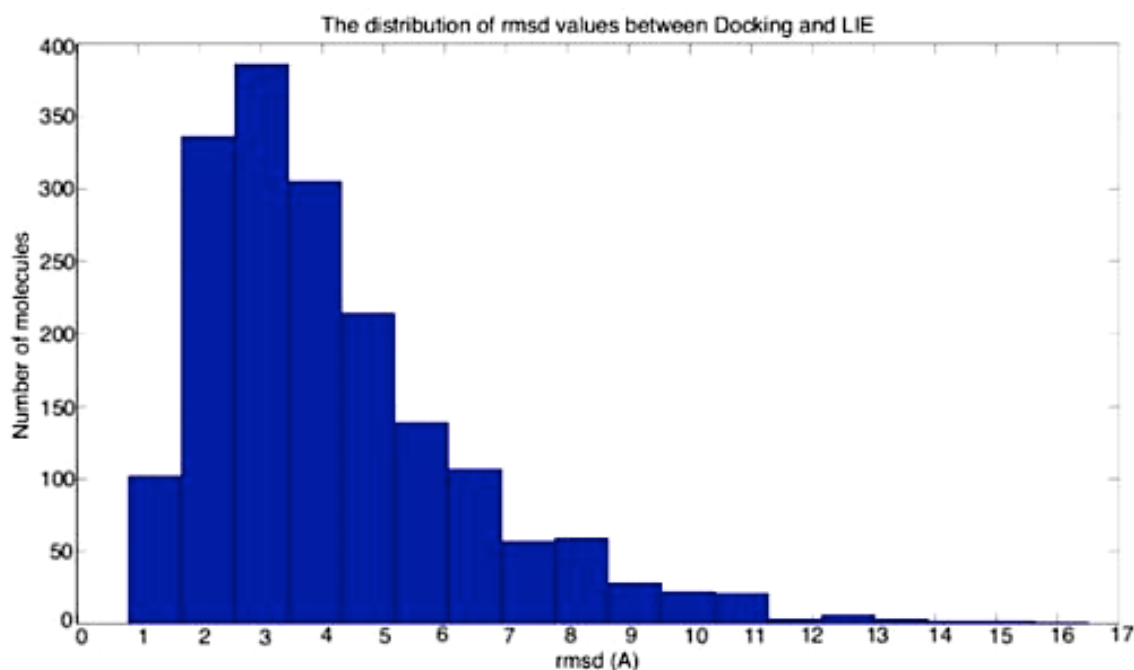
**Figure 56:** Scatter plots showing the correlation of binding free energies (kcal/mol) calculated with the LIE method with normalized consensus scores and individual scoring functions for **dock3**.

Comparing between docking experiments **dock1** and **dock2** means comparing the two different pharmacophore filters used in this study. Even though the molecules that passed the pharmacophore filter **pharma2** were more than twice the number of molecules that passed **pharma1**, the docking experiments **dock1** and **dock2** didn't yield very different results. There were 150 intersecting molecules in the selected sets of **dock1** and **dock2**.

On the other hand, comparing **dock1** and **dock3** enables the comparison of docking programs AutoDock Vina and Surflex-Dock since both experiments employed the same pharmacophore filter, **pharma1**. Since correlation values between the *NCS* and the individual functions are higher in the case of **dock3**, it can be concluded that Surflex-Dock was able to create binding modes that were overall more favorable by the individual scoring functions.

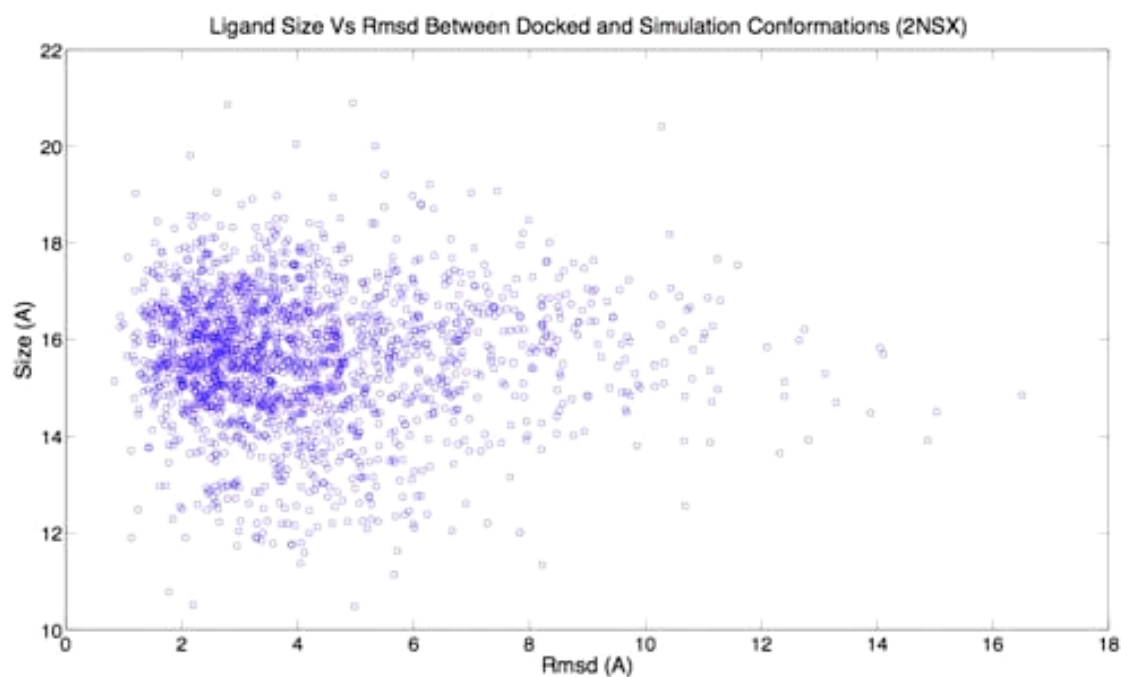


The rmsd comparison between the binding modes output by the docking programs and by the molecular dynamics simulations done in the solvated spherical system centered on the protein binding site, showed that around 60% of the molecules have similar binding modes (Figure 57). `g_rms` command of GROMACS<sup>188</sup> was used for rmsd calculations. The rmsd values between the docked conformations and final simulation conformations show that most of the ligands are 1-4 Å apart. Given the difference of the allowed regions, these rmsd values are small enough to conclude that docking and simulations created similar binding modes for over 1100 ligands. On the other hand, the rest of the ligands show larger displacements from their docked binding modes after molecular dynamics simulations. One reason for these large values of rmsd can be again the difference of the allowed regions in docking and in simulations. In molecular dynamics simulations, the ligands are free to leave the binding site; however in docking studies they are not. The other reason may be the existence of water molecules. In docking studies, there was no water to interfere with the binding; however LIE simulations took place in aqueous environment. There might be water molecules making bonds with the binding site residues, thus causing the ligands shift their locations.



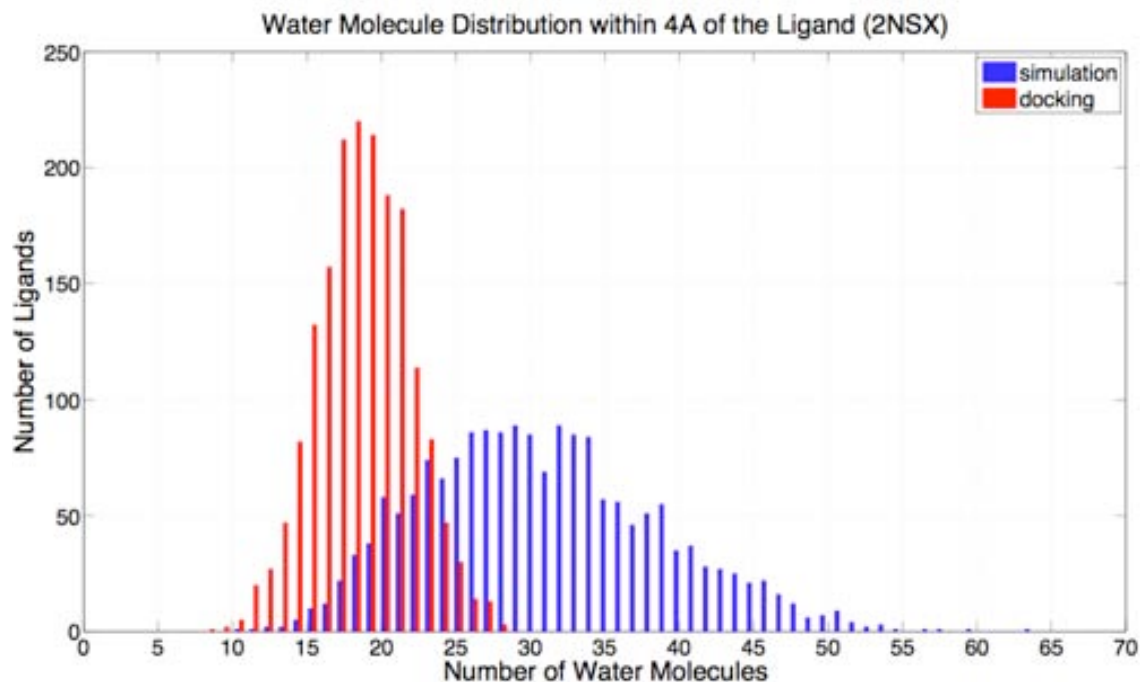
**Figure 57:** Rmsd differences between the docking results and LIE simulation results for GCase.

The plot of the ligand size versus rmsd difference of docked and simulation binding modes for human GCase also showed that a visible pattern was not existent, as in the cases of human T-protein and human bleomycin hydrolase (Figure 58). However, almost 90% of the ligands showed small rmsd differences independent of ligand size. This shows that most of the molecules could be similarly bound either by docking or by simulation due to the well-defined pocket-like binding site of GCase.



**Figure 58:** Rmsd differences between the docking ligand configurations and LIE simulation results (configuration of last time frame) versus the ligand size for GCase.

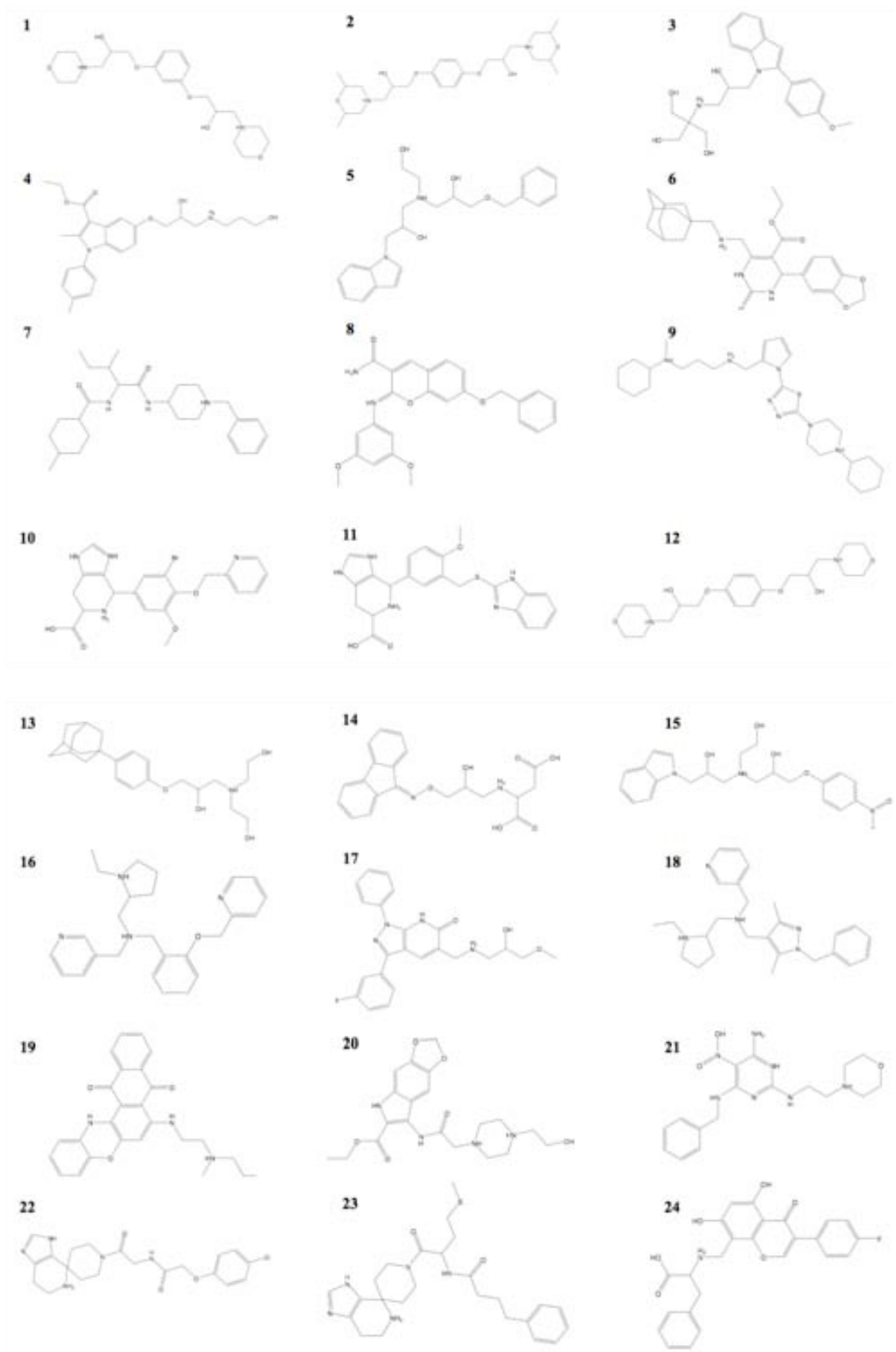
The comparison of number of water molecules around 4 Å of the ligand at the beginning (solvated docked complexes) and end of the LIE simulations shows that the degree of solvation of the ligand increases during the simulation, as observed in the two previous cases (T-protein and hBH). (Figure 59).

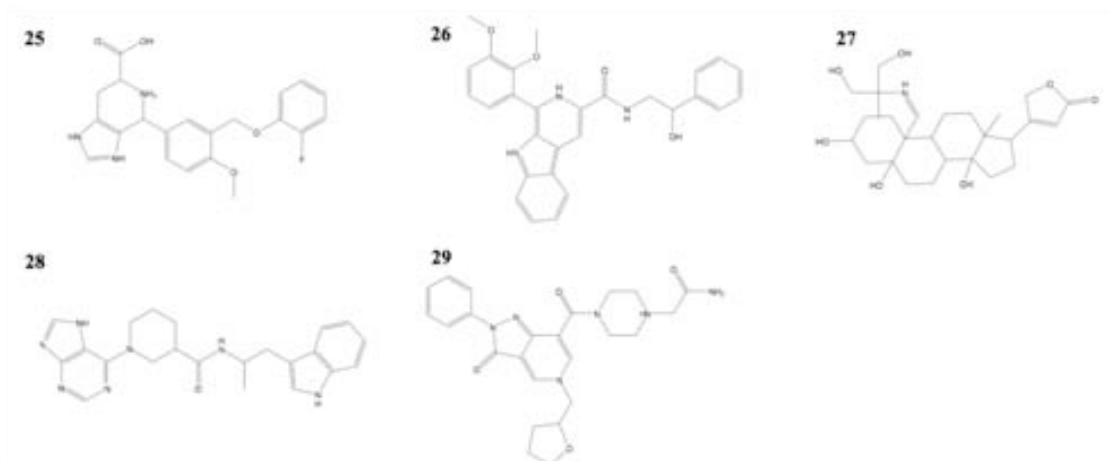


**Figure 59:** The distribution of water molecules within 4 Å of the ligands at the end of docking experiments and LIE simulations for Gcase.

To choose for the candidate molecules that would continue to experimental testing, we applied three criteria and the molecules that fulfill at least two of these criteria were chosen for experimental testing. The first criterion was passing the visual inspection step. All 1800 conformations were inspected visually using MOE<sup>32</sup>. At the end of visual inspection, 100 molecules were chosen, regarding their abilities to make hydrogen bonds with the binding site residues, their exposure to the solvent and the existence of  $\pi$ - $\pi$  or cation- $\pi$  interactions with the protein. The second determining factor was to be in the intersection set of molecules chosen in the top 600 by all three docking experiments. There were 28 molecules (Figure 53) in total that fulfilled this requirement. Finally, the last test was to be in the top 100 after ranking according to binding free energies calculated with LIE. For this last test, repeating conformations of the same molecule were not treated as separate cases, the conformation with the lowest binding free energy was chosen, and thus 100 different molecules were taken.

To define the list of molecules that would be tested experimentally, we selected 22 molecules that fulfilled at least two of the three requirements. In addition to this, to increase the number of molecules, 7 more molecules ranking in the top 30 according to LIE results, but failing the other two criteria were also added. Therefore, 29 molecules in total were selected for experimental testing (Figure 60). While selecting the molecules, we also considered the rmsd values between docked and simulation conformations, and tried to choose molecules that show similar binding modes in both docking and simulation results (Table 21).





**Figure 60:** Molecules chosen for experimental testi

Molecule Number	Selection Criteria	Rmsd between docking and LIE (Å)
1	LIE, Visual Inspection	1.35
2	LIE*	2.34
3	LIE, Visual Inspection	5.24
4	LIE, Visual Inspection	4.73
5	LIE, Visual Inspection	2.43
6	Three Dockings, Visual Inspection	2.15
7	LIE*	3.33
8	LIE, Visual Inspection	2.50
9	LIE*	1.90
10	LIE, Three Dockings, Visual Inspection	2.72
11	LIE, Three Dockings, Visual Inspection	0.93
12	LIE*	2.54
13	LIE*	6.42
14	LIE, Visual Inspection	6.00
15	LIE, Visual Inspection	3.69
16	LIE, Visual Inspection	1.26
17	LIE, Visual Inspection	2.93
18	LIE*	2.56
19	LIE, Three Dockings, Visual Inspection	2.72
20	LIE, Visual Inspection	2.66
21	LIE, Visual Inspection	3.13

22	LIE, Visual Inspection	2.24
23	LIE, Three Dockings	2.87
24	LIE*	1.95
25	LIE, Three Dockings, Visual Inspection	1.59
26	LIE, Three Dockings	1.17
27	LIE, Visual Inspection	3.74
28	LIE, Visual Inspection	5.38
29	LIE, Visual Inspection	1.73

**Table 21:** The list of selection criteria fulfilled by the chosen molecules. LIE: ranking in the top 100, LIE\*: ranking in the top 30, Visual Inspection: in the top 100 molecules chosen with visual examination, Three Dockings: ranked in the top 600 of all three docking experiments.

The 29 chosen molecules were experimentally tested for activity by the Group of Amadeu Llebaria from the Institute of Advanced Chemistry of Catalonia, Consejo Superior de Investigaciones Científicas (IQAC-CSIC). The determination of the activity of imiglucerase (Cerezyme®, Genzyme), a recombinant analogue of human  $\beta$ -glucocerebrosidase, was performed using a fluorimetric method based on the hydrolysis of 4-methylumbelliferyl- $\beta$ -D-glucopyranoside (4-UMG) to glucose and 4-methylumbelliferone (4-MU, Figure 61), catalyzed by Imiglucerase. Each vial of 200 units contains approximately 5 mg of enzyme.

**Figure 61:** Hydrolysis of the substrate by the imiglucerase activity.

The detection equipment is the SpectraMax M5 (Molecular Devices Corporation) for 96-well plates. The fluorescence measurements were performed at a wavelength of 355 nm excitation and a 460 nm emission.

The compounds were dissolved in DMSO at a concentration of 1 mg / mL. From this the appropriate dilutions were prepared. The test also included the iminosugar N-nonyl-deoxynojirimycin (NN-DNJ) to a concentration of 50  $\mu$ M. This compound has been described as an inhibitor of imiglucerase.

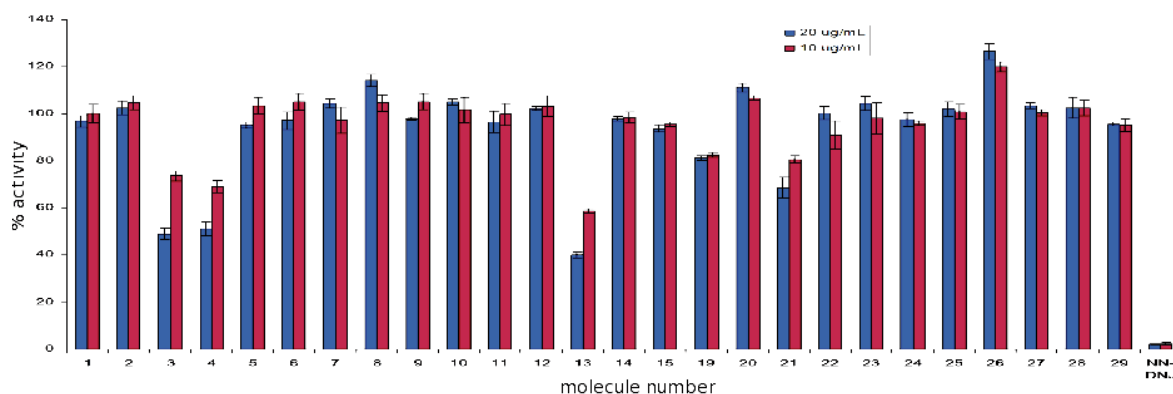
The assays were performed in triplicate. To study a possible dose response, the tests were performed at two concentrations (20 and 10  $\mu$ g / mL). Incubations of 30 minutes were made in the presence and absence of inhibitors at varying concentrations with 25  $\mu$ L of enzyme (0.1 mg protein / mL) in a total volume of 40  $\mu$ L of McIlvaine buffer solution, pH 5.2, 0.1% Triton X-100 (v / v) and 0.2% sodium taurocholate (w / v). Then 60  $\mu$ L of the substrate (4-UMG, 4 mM in McIlvaine buffer solution, pH 5.2) was added and left for reaction for 10 minutes. The incubations were stopped with 150  $\mu$ L of buffer solution glycine / NaOH (100 mM, pH 10.6) (Table 22). Determinations of the 4-MU formed were made at an excitation wavelength of 355

nm and emission 460 nm.

	Control	Inhibitor
Enzyme	25 $\mu$ L	25 $\mu$ L
DMSO	2 $\mu$ L	---
Inhibitor	---	2 $\mu$ L
Buffer solution	13 $\mu$ L	13 $\mu$ L
<b>Preincubation 30 minutes</b>		
Substrate	60 $\mu$ L	60 $\mu$ L
<b>Incubation 10 minutes</b>		
Glycine/NaOH	150 $\mu$ L	150 $\mu$ L

**Table 22:** Volume ratio for activity studies of imiglucrase

Activities for the selected molecules at two concentrations are listed in Table 23 and plotted in Figure 62. Even though the control molecule NN-DNJ performed clearly better than all tested molecules, compounds **3**, **4**, **13**, **19** and **21** affect enzyme activity, while the remaining compounds hardly have any effect. Therefore, it can be concluded that the computational workflow managed to identify five possible hits.



**Figure 62:** Experimentally observed activities of the selected molecules.

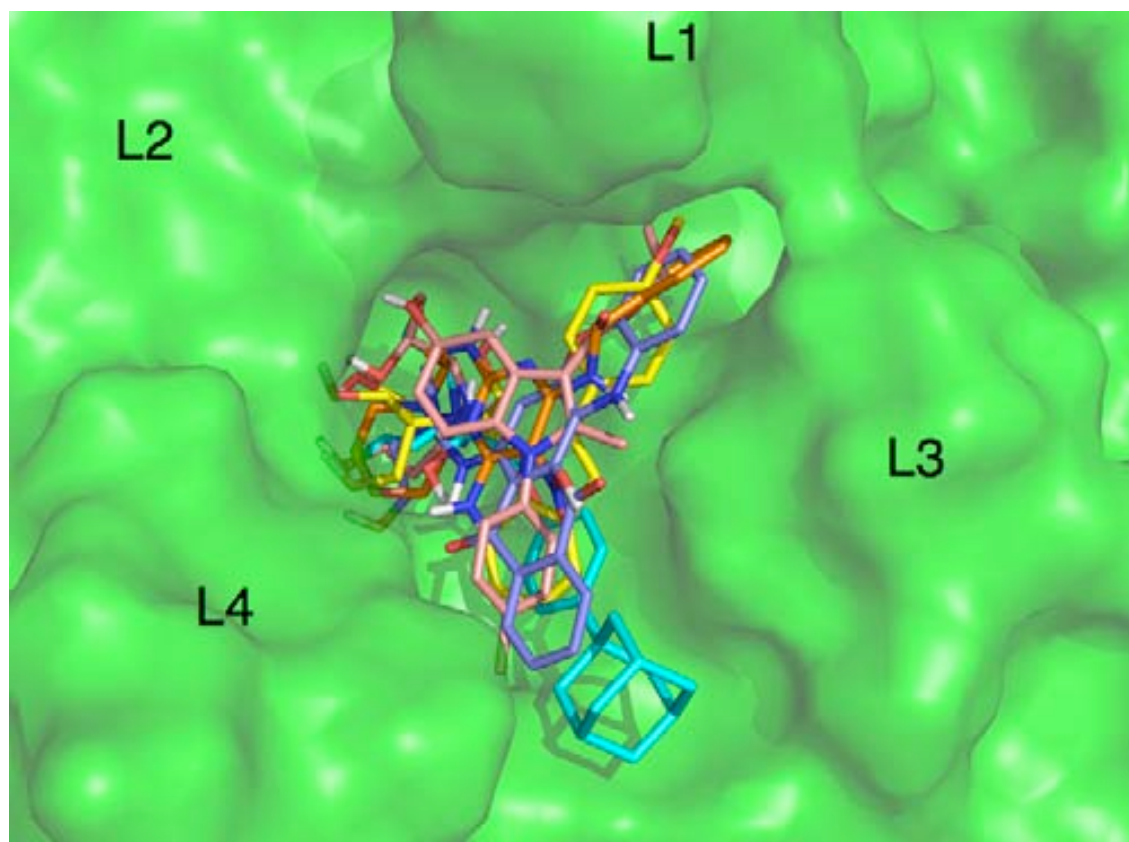
Compound Number	% activity 20 $\mu\text{g/mL}$	% activity 10 $\mu\text{g/mL}$
1	96.8 $\pm$ 2.5	100.1 $\pm$ 3.8
2	102.4 $\pm$ 2.9	104.9 $\pm$ 3.0
3	49.0 $\pm$ 2.4	73.6 $\pm$ 2.1
4	51.3 $\pm$ 2.8	69.1 $\pm$ 2.5
5	95.3 $\pm$ 1.4	103.5 $\pm$ 3.6
6	97.1 $\pm$ 3.7	105.1 $\pm$ 3.5
7	104.4 $\pm$ 2.1	97.1 $\pm$ 5.5
8	114.2 $\pm$ 2.6	104.7 $\pm$ 3.5
9	97.9 $\pm$ 0.6	105.2 $\pm$ 3.3
10	105.1 $\pm$ 1.4	101.8 $\pm$ 5.4
11	96.5 $\pm$ 4.6	99.7 $\pm$ 4.8
12	102.2 $\pm$ 0.9	103.2 $\pm$ 4.5
13	40.2 $\pm$ 1.2	58.8 $\pm$ 1.1
14	97.9 $\pm$ 0.9	98.7 $\pm$ 2.3
15	93.8 $\pm$ 1.3	95.6 $\pm$ 1.1
19	81.3 $\pm$ 1.0	82.5 $\pm$ 0.9
20	111.1 $\pm$ 1.8	106.7 $\pm$ 1.0
21	68.7 $\pm$ 4.3	80.8 $\pm$ 1.3
22	100.3 $\pm$ 2.7	91.1 $\pm$ 5.8
23	104.4 $\pm$ 3.0	98.1 $\pm$ 6.8
24	97.5 $\pm$ 3.0	95.9 $\pm$ 0.9



25	102.0 ± 3.0	100.9 ± 3.0
26	126.5 ± 3.6	120.0 ± 2.0
27	103.5 ± 1.3	100.6 ± 1.2
28	102.6 ± 4.5	102.5 ± 3.2
29	95.8 ± 0.8	95.1 ± 2.7
NN-DNJ (50 μM)	2.2 ± 0.1	2.4 ± 0.4

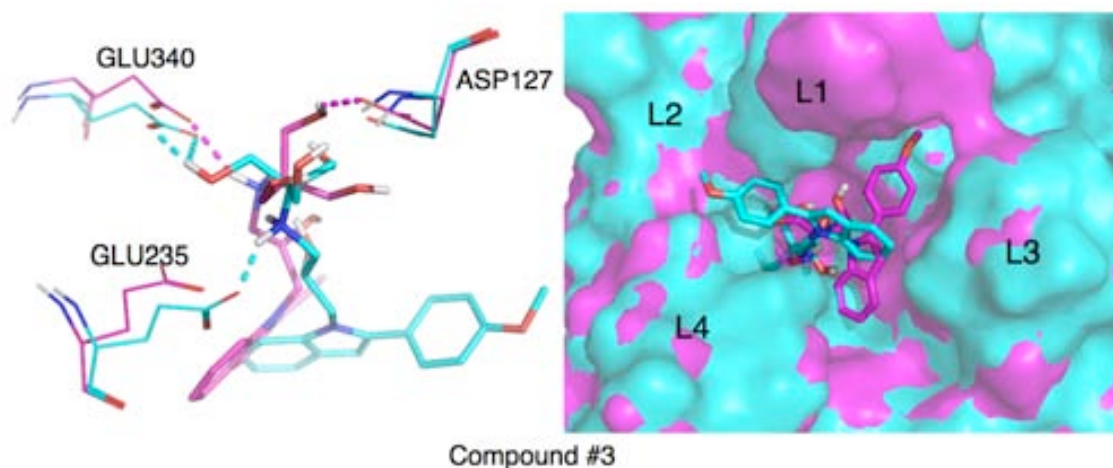
**Table 23:** Experimentally observed activities of the selected molecules.

Docked conformations of the five hit molecules show that while most of the hydrogen bonding interactions occur with binding site residues (Asp127, Trp179, Glu235, Glu340, Trp381 and Asn396), Compounds **3**, **4**, **19** and **21** also occupy the hydrophobic groove between loop L1 (Phe347 and Trp348) and loop L3 (Trp312, Leu314 and Phe316). This suggests that Compounds **3**, **4**, **19** and **21** can establish hydrophobic interactions with the side chains of loops L1 and L3 (Figure 63).



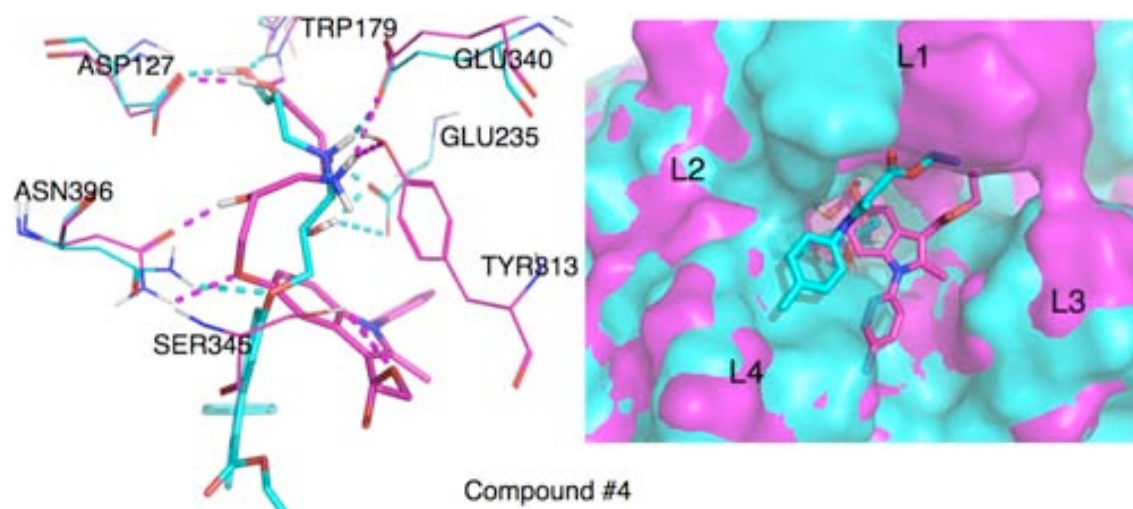
**Figure 63:** Docked binding modes of five hits. Yellow molecule: Compound 3, pink molecule: Compound 4, cyan molecule: Compound 13, purple molecule: Compound 19 and orange molecule: Compound 21. L1-L4 denote the four loops that form the entrance to the binding site. Non-polar hydrogens are not shown for clarity.

When docked and simulation poses of Compound 3 are examined, it can be seen that both poses make hydrogen bonds with the binding site residues. While the docking pose binds to residues Asp127 and Glu340, the simulation pose makes hydrogen bonds to residues Glu235 and Glu340 (left side of Figure 64). Even though hydroxyl moieties of the two poses of the ligand interact with the same residues, the ring parts are positioned quite differently (right side of Figure 64). While the docked conformation lies in the valley between loops L1 and L3, the simulation conformation is positioned between loops L2 and L4. This difference in positioning also explains the large value of rmsd difference (5.24 Å) between docked and simulation conformations of Compound 3.

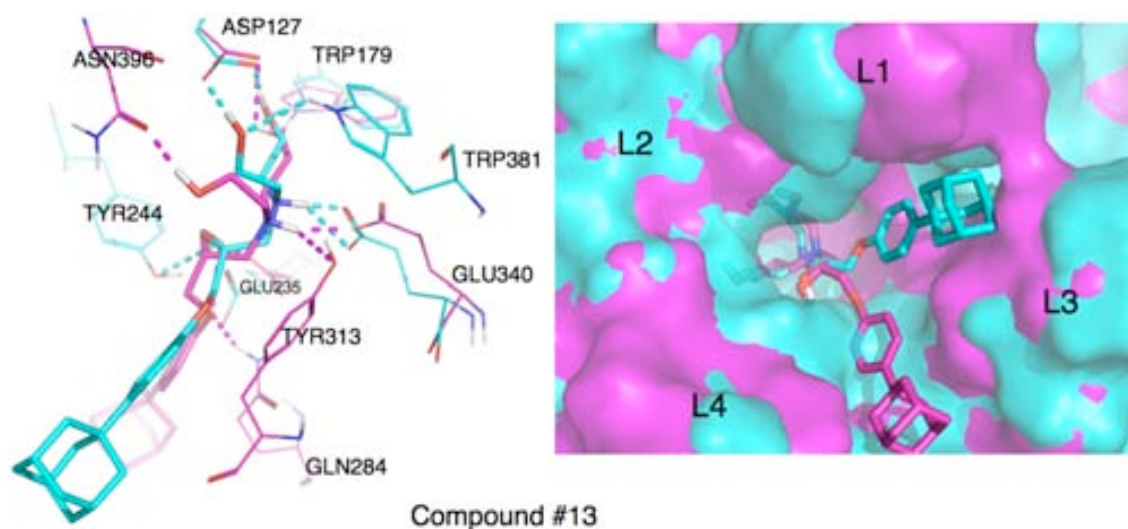


**Figure 64:** Docked and simulation binding modes of Compound 3. The magenta molecule represents the docked conformation of Compound 3, magenta residues are the residues of GCCase structure used in docking and magenta dashed lines are the hydrogen bonds between the docked ligand and the target. The cyan molecule corresponds to the last configuration of Compound 3 from LIE simulations, cyan residues are the last configuration of GCCase from LIE simulations and cyan dashed lines are the hydrogen bonds between Compound 3 and GCCase. Non-polar hydrogens are not shown for clarity.

Docked conformation of Compound 4 makes hydrogen bonds to binding site residues Asp127, Trp179, Glu340 and Asn396, and additionally to residues Tyr313 and Ser345, which are located in the valley between loops L1 and L3 (Figure 65). Simulation binding mode of Compound 4 also hydrogen bonds to the same binding site residues as the docked conformation and additionally to Glu235 (left side of Figure 65). Even though the simulation conformation has a part extending to the area between loops L1 and L3, a hydrogen bonding pattern is not observed, probably due to the movement of L1 as a result molecular dynamics simulations (right side of Figure 65). The 4.23 Å rmsd difference between the docked and simulation binding modes may be a result of the movement of loop L1.

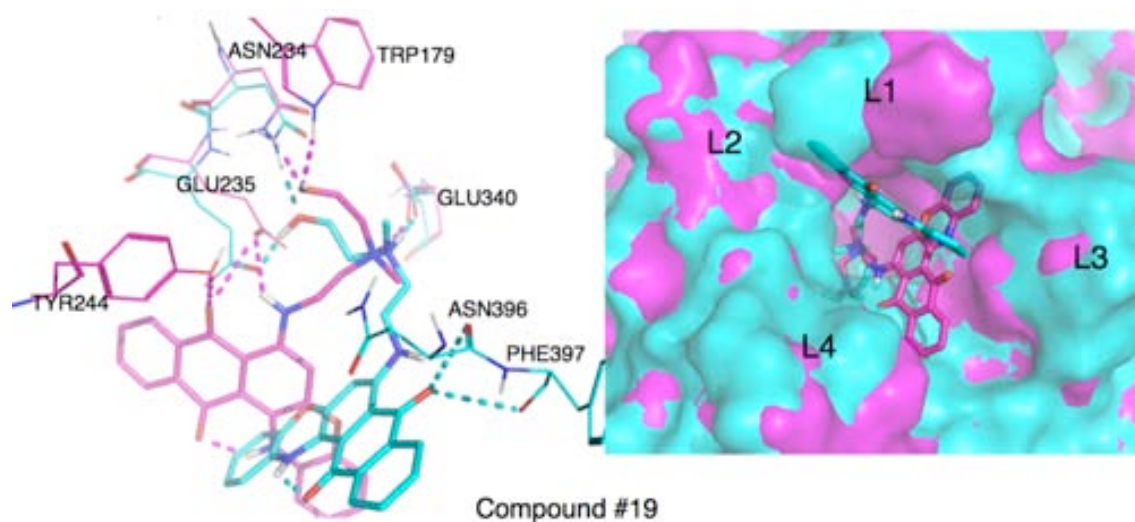


**Figure 65:** Docked and simulation binding modes of Compound 4. Magenta colored binding mode represents the docking result and cyan colored binding mode represents the LIE simulation result. Hydroxyl moieties of both docked and simulation configurations of Compound 13 make hydrogen bonds almost in the same fashion (left side of Figure 66). The docked conformation hydrogen bonds to binding site residues Asp127, Trp179, Glu235, Glu340 and Asn396, while the simulation conformation binds to Asp127, Trp179, Glu235, Glu340 and Trp381. However ring parts of these two conformations are located quite separately in the four-loop area. While the ring part of the docked conformation is located between loops L3 and L4, the ring part of the simulation conformation lies between loops L1 and L3 (right side of Figure 66). If magenta colored protein surface (docking target) is examined in Figure 66, it can be seen that there is a bridge like connection between loops L1 and L3 and it possibly prevents the bulky part of Compound 13 to occupy the valley between L1 and L3. However, the cyan colored simulation conformation shows that L1 and L3 moved away from each other as a result of LIE simulations, enabling the bulky part of Compound 13 to occupy the valley.



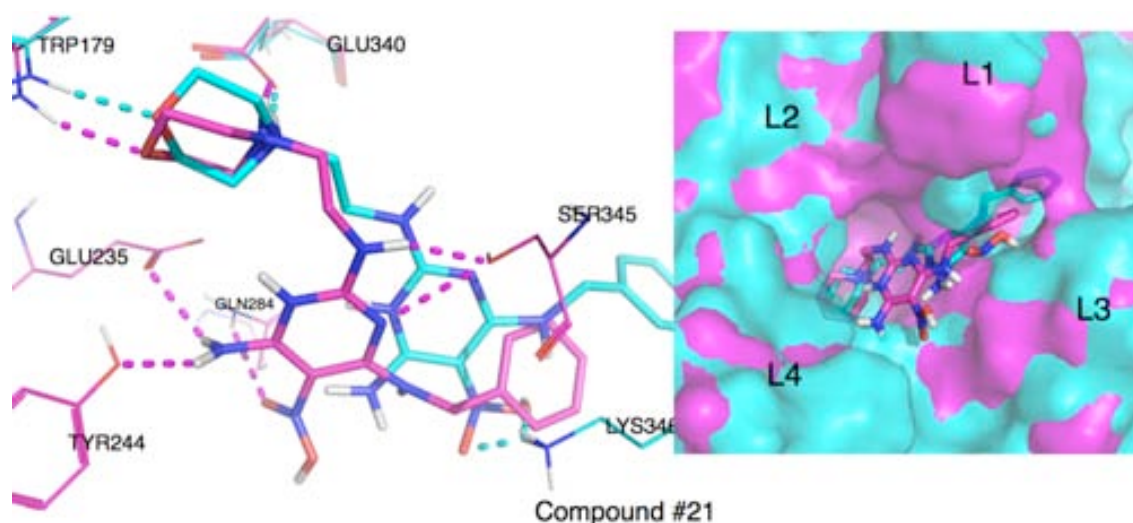
**Figure 66:** Docked and simulation binding modes of Compound 13. Magenta colored binding mode represents the docking result and cyan colored binding mode represents the LIE simulation result.

The docked conformation of Compound 19 makes single hydrogen bonds to binding site residues Trp179 and Glu340, and a double hydrogen bond with Glu235 (left side of Figure 67). The ringed part of docked Compound 19 also interacts with Asn234 and Tyr244 while occupying the area between loops L1 and L3. In the simulation binding mode, Compound 19 interacts with binding site residues Glu235, Glu340 and Asn396 with hydrogen bonds. The ringed part of Compound 19 moves to the area between loops L1 and L2 in the simulation binding mode, probably as a result of the movement of loop L1 (right side of Figure 67).



**Figure 67:** Docked and simulation binding modes of Compound 19. Magenta colored binding mode represents the docking result and cyan colored binding mode represents the LIE simulation result.

In both binding modes, Compound 21 is similarly located in the binding site. The docking binding mode enables hydrogen bonds with residues Trp179, Glu235, Tyr244, Gln284, Glu340 and Ser 345, however the hydrogen bonding pattern in the simulation binding mode is not as extensive (left side of Figure 68). Compound 21 makes hydrogen bonds only with residues Trp179, Glu340 and Lys346 in the simulation binding mode. This is probably due to loop L1 moving away from loop L3, thus expanding the valley in between (right side of Figure 68). As a result of this expansion, Compound 21 has a larger movement area in the simulation binding mode than in the docked binding mode.



**Figure 68:** Docked and simulation binding modes of Compound 21. Magenta colored binding mode represents the docking result and cyan colored binding mode represents the LIE simulation result.

Table 24 shows the hydrogen bonding of the five hit compounds with the GCCase made by both docked and simulation binding modes.

Molecule	Dock	LIE
Compound #3	Asp127 Glu340	Glu235 Glu340 (2)
Compound #4	Asp127 Trp179 Tyr313 Glu340 Ser345 Asn396 (2)	Asp127 Trp179 Glu235 (3) Glu340 Asn396
Compound #13	Asp127	Asp127 (2)

	Trp179 Glu235 Gln284 Tyr313 Glu340 Asn396	Trp179 Glu235 Tyr244 Glu340 Trp381
Compound #19	Trp179 Asn234 Glu235 (2) Tyr244 Glu340	Asn234 Glu235 Glu340 Asn396 Phe397
Compound #21	Trp179 Glu235 Tyr244 Gln284 Glu340 Ser345 (2)	Trp179 Glu340 Lys346

**Table 24:** Hydrogen bonds made between the hit compounds and GCaase. First column is the compound number, second column is the hydrogen bonds made by the docked binding mode and second column is the hydrogen bonds made by the simulation binding mode. The numbers in parenthesis denote the number of hydrogen bonds made with the corresponding residue (if more than 1 bond is made).

#### 4.4. Conclusion

In this study, we proposed a virtual screening procedure that combines pharmacophore design, high-throughput molecular docking, consensus scoring and evaluation of binding free energy by the LIE method. Two large libraries of small molecules have been screened to find potential active binders to two proteins, human T-protein and human bleomycin hydrolase, that are suggested to take part in Alzheimer's Disease, and human acid beta-glucosidase, which is a key protein involved in Gaucher's Disease. Even though human T-protein and human bleomycin hydrolase were not primarily discovered for their involvement in Alzheimer's Disease, their function may be also important in the pathological pathway of the disease and, therefore, they have been proposed as drug targets. For the screening experiment, pharmacophore filtering and molecular docking was employed to reduce the library size and to find possible hits. The pharmacophore filters used did not contain many different features to enable hit variety. The flexible ligand - rigid protein approximation was used for the docking experiments in this study for efficiency. However, this approximation restricts allowed poses for molecules and may cause primarily false negatives but also false positives. The evaluation of binding modes of molecules docked in the protein binding site was done with a modified rank-by-number consensus scoring method. Consensus scoring was used to find out molecules that scored well with all five scoring functions. For each target, the compounds with the best normalized consensus scores were subjected to two molecular dynamics simulations for binding free energy estimation by LIE. LIE based methods are known to perform better than existing scoring functions, mainly because the proteins are not rigid as in docking and the simulations take place in a solvated environment. Therefore,

flexibility to both ligand and protein was introduced during the simulations. Predicted binding free energies were not significantly correlated with either the individual scoring functions or the normalized consensus score for the three screening experiments.

In the continuously developing world of computer-aided drug design, hybrid approaches are needed to compensate for the weaknesses of individual standard methodologies. Increasing computation power enables exploration of rigorous but expensive methods, and to our knowledge, this study is the first application of standard LIE at large scale (around 11800 molecules in total for the three targets). Fully automated treatment of small molecules for different applications makes the workflow explained in this study a very versatile approach for virtual screening of different targets.

#### ***4.5. Future Work***

Candidate ligands need to be chosen for experimental work for the T-protein to see whether there are actives and whether there is correlation between experimental binding affinities and predicted binding free energies by LIE. Additional studies may be done to improve the consensus scoring algorithm, i.e. excluding a scoring function that doesn't correlate with the rest of the scoring functions, or selecting the scoring functions according to the target protein.

# *Chapter 5*

## *Results*





## 5. Conclusions

The main conclusions of the work presented in this thesis can be summarized as follows:

- The design of a drug discovery project requires an integrated approach. For this purpose, we designed an automated computational workflow combining different techniques used in structure-based drug design.
- The computational workflow described here is a result of evolving procedures that started with docking parameter assignment according to the properties of the ligands, and then continued with the three-step docking approach. With the addition of pharmacophore filtering to the three step docking approach, the docking step got more efficient, leading to the hybrid approach. By adding binding free energy predictions with the LIE method, the workflow reached to its final version.
- The computational workflow brings together pharmacophore modeling, high-throughput molecular docking, consensus scoring and binding free energy calculations for small molecules via molecular dynamics simulations.
- The computational workflow is an almost fully automated procedure that only requires manual adjustment at the points of pharmacophore filter creation, grid center definition for docking and charge assignment of the target protein for LIE simulations.
- The workflow is applicable to any protein with a three dimensional structure and a known binding site, either experimentally resolved or computationally predicted. The existence of known inhibitors is not a requirement, however information from the known inhibitor can also be integrated to the workflow if available.
- At the time of writing this thesis, this study represents the first large-scale application (in the thousands of molecules per target) of molecular-dynamics simulation based ligand-binding free energy prediction with the linear interaction energy (LIE) method.



## ***Bibliography***

1. Young, D.C. *Computational Drug Design: A Guide for Computational and Medicinal Chemists*. (Wiley-Interscience: 2009).
2. Gohlke, H. & Klebe, G. Approaches to the description and prediction of the binding affinity of small-molecule ligands to macromolecular receptors. *Angew. Chem. Int. Ed.* **41**, 2644-2676 (2002).
3. Macarron, R. et al. Impact of high-throughput screening in biomedical research. *Nat Rev Drug Discov* **10**, 188-195 (2011).
4. Hertzberg, R.P. & Pope, A.J. High-throughput screening: New technology for the 21st Century. *Current Opinion in Chemical Biology* **4**, 445-451 (2000).
5. Beydon, M.-H., Fournier, A., Drugeault, L. & Becquart, J. Microbiological high throughput screening: An opportunity for the lead discovery process. *Journal of Biomolecular Screening* **5**, 13 -21 (2000).
6. Walters, W.P., Stahl, M.T. & Murcko, M.A. Virtual screening--An overview. *Drug Discovery Today* **3**, 160-178 (1998).
7. Shoichet, B.K. Virtual screening of chemical libraries. *Nature* **432**, 862-865 (2004).
8. Clark, D.E. & Pickett, S.D. Computational methods for the prediction of 'drug-likeness'. *Drug Discovery Today* **5**, 49-58 (2000).
9. Greer, J., Erickson, J.W., Baldwin, J.J. & Varney, M.D. Application of the three-dimensional structures of protein target molecules in structure-based drug design. *Journal of Medicinal Chemistry* **37**, 1035-1054 (1994).
10. Gambacorti-Passerini, C. Part i: Milestones in personalised medicine--imatinib. *Lancet Oncol* **9**, 600 (2008).
11. Druker, B.J. & Lydon, N.B. Lessons learned from the development of an abl tyrosine kinase inhibitor for chronic myelogenous leukemia. *J. Clin. Invest.* **105**, 3-7 (2000).
12. Tagamet: A medicine that changed people's lives. at <http://acswebcontent.acs.org/landmarks/tagamet/tagamet.html>
13. Lemmer, B. The sleep-wake cycle and sleeping pills. *Physiol. Behav* **90**, 285-293 (2007).
14. Meindl, P., Bodo, G., Palese, P., Schulman, J. & Tuppy, H. Inhibition of neuraminidase activity by derivatives of 2-deoxy-2,3-dehydro-n-acetylneuraminic acid. *Virology* **58**, 457-463 (1974).
15. Savarino, A. A historical sketch of the discovery and development of hiv-1 integrase inhibitors. *Expert Opin Investig Drugs* **15**, 1507-1522 (2006).

16. Activity of enfuvirtide-based antiretroviral therapy: Abstract and introduction. at <<http://www.medscape.com/viewarticle/451627>>
17. Oprea, T.I. & Matter, H. Integrating virtual screening in lead discovery. *Current Opinion in Chemical Biology* **8**, 349-358 (2004).
18. Lyne, P.D. Structure-based virtual screening: An overview. *Drug Discovery Today* **7**, 1047-1055 (2002).
19. Klopmand, G. Concepts and applications of molecular similarity. *J. Comput. Chem.* **13**, 539-540 (1992).
20. Willett, P., Barnard, J.M. & Downs, G.M. Chemical similarity searching. *Journal of Chemical Information and Computer Sciences* **38**, 983-996 (1998).
21. Wolber, G. & Langer, T. Ligandscout: 3-D pharmacophores derived from protein-bound ligands and their use as virtual screening filters. *Journal of Chemical Information and Modeling* **45**, 160-169 (2005).
22. Hansch, C., Leo, A. & Hoekman, D.H. *Exploring QSAR.: Fundamentals and applications in chemistry and biology*. (An American Chemical Society Publication: 1995).
23. Langer, T., Hoffmann, R.D., Bachmair, F. & Begle, S. Chemical function based pharmacophore models as suitable filters for virtual 3D-database screening. *Journal of Molecular Structure: THEOCHEM* **503**, 59-72 (2000).
24. Bohm, H.-J., Schneider, G., Mannhold, R., Kubinyi, H. & Folkers, G. *Protein-Ligand Interactions: From Molecular Recognition to Drug Design*. (Wiley-VCH: 2003).
25. Kubinyi, H. Qsar and 3D Qsar in drug design part 1: Methodology. *Drug discovery today* **2**, 457-467 (1997).
26. Kuntz, I.D. Structure-based strategies for drug design and discovery. *Science* **257**, 1078 - 1082 (1992).
27. Brady, G.P. & Stouten, P.F.W. Fast prediction and visualization of protein binding pockets with pass. *Journal of Computer-Aided Molecular Design* **14**, 383-401 (2000).
28. Weisel, M., Proschak, E. & Schneider, G. Pocketpicker: Analysis of ligand binding-sites with shape descriptors. *Chemistry Central Journal* **1**, 7 (2007).
29. Hendlich, M., Rippmann, F. & Barnickel, G. LIGSITE: Automatic and efficient detection of potential small molecule-binding sites in proteins. *Journal of Molecular Graphics and Modelling* **15**, 359-363 (1997).
30. Huang, D. & Caflisch, A. Library screening by fragment-based docking. *J. Mol. Recognit* **23**, 183-193 (2010).

31. *Catalyst*. (Accelrys, Inc.: 10188 Telesis Court, Suite 100 San Diego, CA 92121 USA.).
32. *MOE*. (Chemical Computing Group: 1010 Sherbrooke St. W, Suite 910 Montreal, Quebec, Canada H3A 2R7).
33. *Phase*. (Schrodinger, LLC: 120 West 45th Street, 32nd Floor, New York, NY 10036-4041, USA, ).
34. *SYBYL-X*. (Tripos International: 1699 South Hanley Road, St. Louis, MO 63144-2319, USA, ).
35. *LigandScout*. (Inte:Ligand: Clemens Maria Hofbauer-G. 6 A-2344 Maria Enzersdorf Austria, ).
36. Hein, M., Zilian, D. & Sotriffer, C.A. Docking compared to 3D-pharmacophores: the scoring function challenge. *Drug Discovery Today: Technologies* **7**, e229-e236 (2010).
37. Kitchen, D.B., Decornez, H., Furr, J.R. & Bajorath, J. Docking and scoring in virtual screening for drug discovery: methods and applications. *Nat Rev Drug Discov* **3**, 935-949 (2004).
38. Halperin, I., Ma, B., Wolfson, H. & Nussinov, R. Principles of docking: An overview of search algorithms and a guide to scoring functions. *Proteins* **47**, 409-443 (2002).
39. Brooijmans, N. & Kuntz, I.D. Molecular recognition and docking algorithms. *Annu. Rev. Biophys. Biomol. Struct.* **32**, 335-373 (2011).
40. Morris, G.M. et al. Automated docking using a Lamarckian genetic algorithm and an empirical binding free energy function. *J. Comput. Chem.* **19**, 1639-1662 (1998).
41. Verdonk, M.L., Cole, J.C., Hartshorn, M.J., Murray, C.W. & Taylor, R.D. Improved protein–ligand docking using GOLD. *Proteins* **52**, 609-623 (2003).
42. Leach, A.R. *Molecular modelling: principles and applications*. (Addison-Wesley Longman Ltd: 2001).
43. Caflisch, A. Computational combinatorial ligand design: Application to human thrombin. *J Computer-Aided Mol Des* **10**, 372-396 (1996).
44. Schneider, G., Hartenfeller, M. & Proschak, E. De novo drug design. 165-185doi:10.1002/9780470584170.ch6
45. Rarey, M., Kramer, B., Lengauer, T. & Klebe, G. A fast flexible docking method using an incremental construction algorithm. *Journal of Molecular Biology* **261**, 470–489 (1996).

46. Ewing, T.J.A., Makino, S., Skillman, A.G. & Kuntz, I.D. DOCK 4.0: Search strategies for automated molecular docking of flexible molecule databases. *Journal of Computer-Aided Molecular Design* **15**, 411-428 (2001).
47. Kramer, B., Rarey, M. & Lengauer, T. Evaluation of the FLEXX incremental construction algorithm for protein–ligand docking. *Proteins* **37**, 228-241 (1999).
48. Friesner, R.A. et al. Glide: A new approach for rapid, accurate docking and scoring. 1. Method and assessment of docking accuracy. *Journal of Medicinal Chemistry* **47**, 1739-1749 (2004).
49. Trott, O. & Olson, A.J. AutoDock Vina: Improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *J. Comput. Chem.* **31**, 455-461 (2010).
50. Kwan, J.C., Eksioglu, E.A., Liu, C., Paul, V.J. & Luesch, H. Grassystatins A–C from marine cyanobacteria, potent cathepsin E inhibitors that reduce antigen presentation. *Journal of Medicinal Chemistry* **52**, 5732-5747 (2009).
51. Takatsuka, Y., Chen, C. & Nikaido, H. Mechanism of recognition of compounds of diverse structures by the multidrug efflux pump AcrB of *Escherichia coli*. *Proceedings of the National Academy of Sciences* **107**, 6559 -6565 (2010).
52. Eldstrom, J. & Fedida, D. Modeling of high-affinity binding of the novel atrial anti-arrhythmic agent, vernakalant, to Kv1.5 channels. *Journal of Molecular Graphics and Modelling* **28**, 226-235 (2009).
53. Kellenberger, E., Rodrigo, J., Muller, P. & Rognan, D. Comparative evaluation of eight docking tools for docking and virtual screening accuracy. *Proteins* **57**, 225-242 (2004).
54. Bissantz, C., Folkers, G. & Rognan, D. Protein-based virtual screening of chemical databases. 1. Evaluation of different docking/scoring combinations. *Journal of Medicinal Chemistry* **43**, 4759-4767 (2000).
55. Warren, G.L. et al. A critical assessment of docking programs and scoring functions. *Journal of Medicinal Chemistry* **49**, 5912-5931 (2006).
56. Böhm, H.-J. & Klebe, G. What can we learn from molecular recognition in protein–ligand complexes for the design of new drugs? *Angew. Chem. Int. Ed. Engl.* **35**, 2588-2614 (1996).
57. *IUPAC Compendium of Chemical Terminology*. (IUPAC: Research Triangle Park, NC, 2009).at <<http://goldbook.iupac.org>>
58. Chandler, D. Interfaces and the driving force of hydrophobic assembly. *Nature* **437**, 640-647 (2005).

59. Dougherty, D.A. Cation- $\pi$  interactions in chemistry and biology: A new view of benzene, Phe, Tyr, and Trp. *Science* **271**, 163 -168 (1996).
60. Mecozzi, S., West, A.P. & Dougherty, D.A. Cation- $\pi$  interactions in simple aromatics: electrostatics provide a predictive tool. *Journal of the American Chemical Society* **118**, 2307-2308 (1996).
61. Berman, H.M. Hydrogen bonding in biological structures. G. A. Jeffrey and W. Saenger. *Biophys J* **64**, 1976-1976 (1993).
62. Liljefors, T., Jørgensen, F.S. & Krogsgaard-Larsen, P. *Rational molecular design in drug research: proceedings of a symposium held at the Royal Danish Academy of Sciences and Letters, June 8-12, 1997*. (Munksgaard: 1998).
63. Brown, R.D. & Martin, Y.C. Use of structure-activity data to compare structure-based clustering methods and descriptors for use in compound selection. *Journal of Chemical Information and Computer Sciences* **36**, 572-584 (1996).
64. Humblet, C. & Dunbar Jr, J.B. . 3D Database searching and docking strategies. *Annual Reports in Medicinal Chemistry* **28**, 275-284 (1993).
65. Lemmen, C. & Lengauer, T. Computational methods for the structural alignment of molecules. *Journal of Computer-Aided Molecular Design* **14**, 215-232 (2000).
66. Brandsdal, B.O. et al. Free energy calculations and ligand binding. *Protein Simulations* **66**, 123-158 (2003).
67. Kirkwood, J.G. Statistical mechanics of fluid mixtures. *The Journal of Chemical Physics* **3**, 300-313 (1935).
68. Zwanzig, R.W. High-temperature equation of state by a perturbation method. i. nonpolar gases. *J. Chem. Phys.* **22**, 1420 (1954).
69. Åqvist, J., Medina, C. & Samuelsson, J.-E. A new method for predicting binding affinity in computer-aided drug design. *Protein Engineering* **7**, 385 -391 (1994).
70. Hansson, T., Marelus, J. & Åqvist, J. Ligand binding affinity prediction by linear interaction energy methods. *Journal of Computer-Aided Molecular Design* **12**, 27-35 (1998).
71. Srinivasan, J., Cheatham, T.E., Cieplak, P., Kollman, P.A. & Case, D.A. Continuum solvent studies of the stability of DNA, RNA, and phosphoramidate-DNA helices. *Journal of the American Chemical Society* **120**, 9401-9409 (1998).
72. Kollman, P.A. et al. Calculating structures and free energies of complex molecules: combining molecular mechanics and continuum models. *Accounts of Chemical Research* **33**, 889-897 (2000).



73. Englebienne, P. & Moitessier, N. Docking ligands into flexible and solvated macromolecules. 4. are popular scoring functions accurate for this class of proteins? *Journal of Chemical Information and Modeling* **49**, 1568-1580 (2009).
74. Bohm, H.-J. LUDI: Rule-based automatic design of new substituents for enzyme inhibitor leads. *J Computer-Aided Mol Des* **6**, 593-606 (1992).
75. Eldridge, M.D., Murray, C.W., Auton, T.R., Paolini, G.V. & Mee, R.P. Empirical scoring functions: I. The development of a fast empirical scoring function to estimate the binding affinity of ligands in receptor complexes. *Journal of Computer-Aided Molecular Design* **11**, 425-445 (1997).
76. Böhm, H.-J. The development of a simple empirical scoring function to estimate the binding constant for a protein-ligand complex of known three-dimensional structure. *Journal of Computer-Aided Molecular Design* **8**, 243-256 (1994).
77. Böhm, H.-J. Prediction of binding constants of protein ligands: A fast method for the prioritization of hits obtained from de novo design or 3D database search programs. *Journal of Computer-Aided Molecular Design* **12**, 309 (1998).
78. Huey, R., Morris, G.M., Olson, A.J. & Goodsell, D.S. A semiempirical free energy force field with charge-based desolvation. *J. Comput. Chem.* **28**, 1145-1152 (2007).
79. Berman, H.M. et al. The Protein Data Bank. *Nucleic Acids Research* **28**, 235 -242 (2000).
80. Muegge, I. A knowledge-based scoring function for protein-ligand interactions: Probing the reference state. *Perspectives in Drug Discovery and Design* **20**, 99-114 (2000).
81. Muegge, I. & Martin, Y.C. A general and fast scoring function for protein-ligand interactions: a simplified potential approach. *Journal of Medicinal Chemistry* **42**, 791-804 (1999).
82. Muegge, I. PMF scoring revisited. *Journal of Medicinal Chemistry* **49**, 5895-5902 (2006).
83. Gohlke, H., Hendlich, M. & Klebe, G. Knowledge-based scoring function to predict protein-ligand interactions. *Journal of Molecular Biology* **295**, 337-356 (2000).
84. Wang, R., Lu, Y. & Wang, S. Comparative evaluation of 11 scoring functions for molecular docking. *Journal of Medicinal Chemistry* **46**, 2287-2303 (2003).
85. Nervall, M., Hanspers, P., Carlsson, J., Boukharta, L. & Åqvist, J. Predicting binding modes from free energy calculations. *Journal of Medicinal Chemistry* **51**, 2657-2667 (2008).
86. Wang, R. & Wang, S. How does consensus scoring work for virtual library screening? an idealized computer experiment. *Journal of Chemical Information and Computer Sciences* **41**, 1422-1426 (2001).

87. Clark, R.D., Strizhev, A., Leonard, J.M., Blake, J.F. & Matthew, J.B. Consensus scoring for ligand/protein interactions. *Journal of Molecular Graphics and Modelling* **20**, 281-295 (2002).
88. Oda, A., Tsuchida, K., Takakura, T., Yamaotsu, N. & Hirono, S. Comparison of consensus scoring strategies for evaluating computational models of protein–ligand complexes. *Journal of Chemical Information and Modeling* **46**, 380-391 (2006).
89. Brown, S.P. & Muchmore, S.W. Large-scale application of high-throughput molecular mechanics with poisson–boltzmann surface area for routine physics-based scoring of protein–ligand complexes. *Journal of Medicinal Chemistry* **52**, 3159-3165 (2009).
90. Kuhn, B., Gerber, P., Schulz-Gasch, T. & Stahl, M. Validation and use of the MM-PBSA approach for drug discovery. *Journal of Medicinal Chemistry* **48**, 4040-4048 (2005).
91. Stjernschantz, E. et al. Are automated molecular dynamics simulations and binding free energy calculations realistic tools in lead optimization? An evaluation of the linear interaction energy (lie) method. *Journal of Chemical Information and Modeling* **46**, 1972-1983 (2006).
92. Carlsson, J., Boukharta, L. & Åqvist, J. Combining docking, molecular dynamics and the linear interaction energy method to predict binding modes and affinities for non-nucleoside inhibitors to hiv-1 reverse transcriptase. *Journal of Medicinal Chemistry* **51**, 2648-2656 (2008).
93. Bjelic, S. et al. Computational inhibitor design against malaria plasmepsins. *Cell. Mol. Life Sci.* **64**, 2285-2305 (2007).
94. Wallin, G., Nervall, M., Carlsson, J. & Åqvist, J. Charges for large scale binding free energy calculations with the linear interaction energy method. *Journal of Chemical Theory and Computation* **5**, 380-395 (2009).
95. Åqvist, J. & Hansson, T. On the validity of electrostatic linear response in polar solvents. *J. Phys. Chem.* **100**, 9512-9521 (2011).
96. Almlöf, M., Brandsdal, B.O. & Aqvist, J. Binding affinity prediction with different force fields: examination of the linear interaction energy method. *J Comput Chem* **25**, 1242-1254 (2004).
97. Griffith, O.W. & Stuehr, D.J. Nitric oxide synthases: properties and catalytic mechanism. *Annu. Rev. Physiol* **57**, 707-736 (1995).
98. Hou, Y.C., Janczuk, A. & Wang, P.G. Current trends in the development of nitric oxide donors. *Curr. Pharm. Des* **5**, 417-441 (1999).
99. Stuehr, D.J. Mammalian nitric oxide synthases. *Biochim. Biophys. Acta* **1411**, 217-230 (1999).

100. Knowles, R.G. & Moncada, S. Nitric oxide synthases in mammals. *Biochem J* **298**, 249-258 (1994).
101. Griffith, O.W. & Kilbourn, R.G. Design of nitric oxide synthase inhibitors and their use to reverse hypotension associated with cancer immunotherapy. *Advances in Enzyme Regulation* **37**, 171-194 (1997).
102. Marsden, P.A. et al. Molecular cloning and characterization of human endothelial nitric oxide synthase. *FEBS Lett* **307**, 287-293 (1992).
103. Yoshimura, M. et al. A missense Glu298Asp variant in the endothelial nitric oxide synthase gene is associated with coronary spasm in the Japanese. *Hum. Genet* **103**, 65-69 (1998).
104. Wang, B. et al. Association between Alzheimer's disease and the NOS3 gene Glu298Asp polymorphism in Chinese. *J. Mol. Neurosci* **34**, 173-176 (2008).
105. Rosenfeld, R.J. et al. Conformational changes in nitric oxide synthases induced by chlorzoxazone and nitroindazoles: crystallographic and computational analyses of inhibitor potency. *Biochemistry* **41**, 13915-13925 (2002).
106. Geoffrey M Cooper The Eukaryotic Cell Cycle. (2000). at <http://www.ncbi.nlm.nih.gov/books/NBK9876/>
107. Nurse, P. Checkpoint pathways come of age. *Cell* **91**, 865-867 (1997).
108. Hanahan, D. & Weinberg, R.A. The hallmarks of cancer. *Cell* **100**, 57-70 (2000).
109. Zhou, B.-B.S. & Elledge, S.J. The DNA damage response: putting checkpoints in perspective. *Nature* **408**, 433-439 (2000).
110. Cuddihy, A.R. & O'Connell, M.J. Cell-cycle responses to DNA damage in G2. *A Survey of Cell Biology* **Volume 222**, 99-140 (2003).
111. Foloppe, N. et al. Structure-based design of novel chk1 inhibitors: insights into hydrogen bonding and protein-ligand affinity. *Journal of Medicinal Chemistry* **48**, 4332-4345 (2005).
112. Lyne, P.D. et al. Identification of compounds with nanomolar binding affinity for checkpoint kinase-1 using knowledge-based virtual screening. *Journal of Medicinal Chemistry* **47**, 1962-1968 (2004).
113. Rhind, N. & Russell, P. Chk1 and Cds1: Linchpins of the DNA damage and replication checkpoint pathways. *Journal of Cell Science* **113**, 3889-3896 (2000).
114. Boddy, M.N., Furnari, B., Mondesert, O. & Russell, P. Replication checkpoint enforced by kinases Cds1 and Chk1. *Science* **280**, 909-912 (1998).

115. Foloppe, N. et al. Identification of chemically diverse Chk1 inhibitors by receptor-based virtual screening. *Bioorg. Med. Chem* **14**, 4792-4802 (2006).
116. Chen, P. et al. Implications for Chk1 Regulation: The 1.7 Å crystal structure of human cell cycle checkpoint kinase Chk1. *Cell* **100**, 681-692 (2000).
117. Foloppe, N. et al. Identification of a buried pocket for potent and selective inhibition of Chk1: Prediction and verification. *Bioorganic & Medicinal Chemistry* **14**, 1792-1804 (2006).
118. *LigPrep*. (Schrodinger, LLC: 120 West 45th Street, 32nd Floor, New York, NY 10036-4041, USA, ).
119. Cornell, W.D. et al. a second generation force field for the simulation of proteins, nucleic acids, and organic molecules. *Journal of the American Chemical Society* **117**, 5179-5197 (1995).
120. Lipinski, C.A., Lombardo, F., Dominy, B.W. & Feeney, P.J. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Advanced Drug Delivery Reviews* **46**, 3-26 (2001).
121. Veber, D.F. et al. Molecular properties that influence the oral bioavailability of drug candidates. *Journal of Medicinal Chemistry* **45**, 2615-2623 (2002).
122. What is the format of a PDBQT file? — AutoDock. at <<http://autodock.scripps.edu/faqs-help/faq/what-is-the-format-of-a-pdbqt-file>>
123. Sanner, M.F. Python: a programming language for software integration and development. *J. Mol. Graphics Mod* **17**, 57-61 (1999).
124. Sample Mol2 File. at <[http://tripos.com/mol2/mol2\\_format3.html](http://tripos.com/mol2/mol2_format3.html)>
125. Gasteiger, J. & Marsili, M. Iterative partial equalization of orbital electronegativity--a rapid access to atomic charges. *Tetrahedron* **36**, 3219-3228 (1980).
126. Zhao, B. et al. Structural basis for Chk1 inhibition by UCN-01. *J. Biol. Chem* **277**, 46609-46615 (2002).
127. PubChem: Integrated platform of small molecules and biological activities. at <<http://www.acscomp.org/Publications/ARCC/volume4/chapter12.html>>
128. Allain, H. et al. Alzheimer's disease: the pharmacological pathway. *Fundamental & Clinical Pharmacology* **17**, 419-428 (2003).
129. Okamura-Ikeda, K. et al. Crystal structure of human t-protein of glycine cleavage system at 2.0 Å resolution and its implication for understanding non-ketotic hyperglycinemia. *Journal of Molecular Biology* **351**, 1146-1159 (2005).

130. O'Farrell, P.A., Gonzalez, F., Zheng, W., Johnston, S.A. & Joshua-Tor, L. Crystal structure of human bleomycin hydrolase, a self-compartmentalizing cysteine protease. *Structure* **7**, 619-627 (1999).
131. Lieberman, R.L. et al. Structure of acid [beta]-glucosidase with pharmacological chaperone provides insight into Gaucher disease. *Nat Chem Biol* **3**, 101-107 (2007).
132. Fujiwara, K., Okamura-Ikeda, K. & Motokawa, Y. Mechanism of the glycine cleavage reaction. Further characterization of the intermediate attached to H-protein and of the reaction catalyzed by T-protein. *Journal of Biological Chemistry* **259**, 10664-10668 (1984).
133. Okamura-Ikeda, K., Fujiwara, K. & Motokawa, Y. Mechanism of the glycine cleavage reaction. Properties of the reverse reaction catalyzed by T-protein. *Journal of Biological Chemistry* **262**, 6746-6749 (1987).
134. Tada, K. & Kure, S. Non-ketotic hyperglycinaemia: Molecular lesion, diagnosis and pathophysiology. *Journal of Inherited Metabolic Disease* **16**, 691-703 (1993).
135. Kure, S. et al. A one-base deletion (183delC) and a missense mutation (D276H) in the T-protein gene from a Japanese family with nonketotic hyperglycinemia. *J Hum Genet* **43**, 135-137 (1998).
136. Kure, S. et al. A missense mutation (His42Arg) in the T-protein gene from a large Israeli-Arab kindred with nonketotic hyperglycinemia. *Human Genetics* **102**, 430-434 (1998).
137. Toone, J.R., Applegarth, D.A., Levy, H.L., Coulter-Mackie, M.B. & Lee, G. Molecular genetic and potential biochemical characteristics of patients with T-protein deficiency as a cause of glycine encephalopathy (NKH). *Molecular Genetics and Metabolism* **79**, 272-280 (2003).
138. Nanao, K. et al. Identification of the mutations in the t-protein gene causing typical and atypical nonketotic hyperglycinemia. *Human Genetics* **93**, 655-658 (1994).
139. Toone, J.R., Applegarth, D.A., Coulter-Mackie, M.B. & James, E.R. Biochemical and molecular investigations of patients with nonketotic hyperglycinemia. *Molecular Genetics and Metabolism* **70**, 116-121 (2000).
140. Seiler, N. Ammonia and Alzheimer's disease. *Neurochemistry International* **41**, 189-207.
141. Johnson, J.W. & Ascher, P. Glycine potentiates the NMDA response in cultured mouse brain neurons. *Nature* **325**, 529-531 (1987).
142. Kemp, J.A. & Leeson, P.D. The glycine site of the NMDA receptor -- five years on. *Trends in Pharmacological Sciences* **14**, 20-25 (1993).

143. Leeson, P.D. & Iversen, L.L. The glycine site on the NMDA Receptor: Structure-activity relationships and therapeutic potential. *Journal of Medicinal Chemistry* **37**, 4053-4067 (1994).
144. D'Souza, D.C., Charney, D. & Krystal, J. Glycine site agonists of the NMDA receptor: A review. *CNS Drug Reviews* **1**, 227-260 (1995).
145. Tsang, S.W.Y. et al. Alterations in NMDA receptor subunit densities and ligand binding to glycine recognition sites are associated with chronic anxiety in Alzheimer's disease. *Neurobiology of Aging* **29**, 1524-1532 (2008).
146. Arlt, M. & Bartoszyk, G. Cleavage system inhibitors as potential antipsychotics. (2002).
147. Douce, R., Bourguignon, J., Neuburger, M. & Rébeillé, F. The glycine decarboxylase system: a fascinating complex. *Trends in Plant Science* **6**, 167-176 (2001).
148. Bromme, D., Rossi, A.B., Smeekens, S.P., Anderson, D.C. & Payan, D.G. Human bleomycin hydrolase: molecular cloning, sequencing, functional expression, and enzymatic characterization. *Biochemistry* **35**, 6706-6714 (1996).
149. Enenkel, C. & Wolf, D.H. BLH1 codes for a yeast thiol aminopeptidase, the equivalent of mammalian bleomycin hydrolase. *Journal of Biological Chemistry* **268**, 7036 -7043 (1993).
150. Berti, P.J. & Storer, A.C. Alignment/phylogeny of the papain superfamily of cysteine proteases. *Journal of Molecular Biology* **246**, 273-283 (1995).
151. Chapman, H.A., Riese, R.J. & Shi, G.-P. Emerging roles for cysteine proteases in human biology. *Annu. Rev. Physiol.* **59**, 63-88 (2011).
152. Sebti, S.M., Mignano, J.E., Jani, J.P., Srimatkandada, S. & Lazo, J.S. Bleomycin hydrolase: molecular cloning, sequencing, and biochemical studies reveal membership in the cysteine proteinase family. *Biochemistry* **28**, 6544-6548 (1989).
153. Sebti, S.M. & Lazo, J.S. Metabolic inactivation of bleomycin analogs by bleomycin hydrolase. *Pharmacology & Therapeutics* **38**, 321-329 (1988).
154. Koldamova, R.P. et al. Human bleomycin hydrolase binds ribosomal proteins†,‡. *Biochemistry* **38**, 7111-7117 (1999).
155. Koldamova, R.P., Lefterov, I.M., DiSabella, M.T. & Lazo, J.S. An Evolutionarily conserved cysteine protease, human bleomycin hydrolase, binds to the human homologue of ubiquitin-conjugating enzyme 9. *Molecular Pharmacology* **54**, 954 -961 (1998).
156. Lefterov, I.M., Koldamova, R.P. & Lazo, J.S. Human bleomycin hydrolase regulates the secretion of amyloid precursor protein. *The FASEB Journal* **14**, 1837 -1847 (2000).
157. Selkoe, D.J. Amyloid  $\beta$ -Protein and the Genetics of Alzheimer's Disease. *Journal of Biological Chemistry* **271**, 18295 -18298 (1996).

158. Namba, Y., Ouchi, Y., Takeda, A., Ueki, A. & Ikeda, K. Bleomycin hydrolase immunoreactivity in senile plaque in the brains of patients with Alzheimer's disease. *Brain Research* **830**, 200-202 (1999).
159. Montoya, S.E. et al. Bleomycin hydrolase is associated with risk of sporadic Alzheimer's disease. *Nat Genet* **18**, 211-212 (1998).
160. Sawkar, A.R., D'Haese, W. & Kelly, J.W. Therapeutic strategies to ameliorate lysosomal storage disorders – a focus on Gaucher disease. *Cell. Mol. Life Sci.* **63**, 1179-1192 (2006).
161. Zhao, H. & Grabowski, G.A. Gaucher disease: perspectives on a prototype lysosomal disease. *Cellular and Molecular Life Sciences (CMLS)* **59**, 694-707 (2002).
162. Jmoudiak, M. & Futerman, A.H. Gaucher disease: pathological mechanisms and modern management. *Br J Haematol* **129**, 178-188 (2005).
163. Futerman, A.H. & van Meer, G. The cell biology of lysosomal storage disorders. *Nat Rev Mol Cell Biol* **5**, 554-565 (2004).
164. Goker-Alpan, O. Optimal therapy in Gaucher disease. *Ther Clin Risk Manag* **6**, 315-323 (2010).
165. Grace, M.E., Newman, K.M., Scheinker, V., Berg-Fussman, A. & Grabowski, G.A. Analysis of human acid beta-glucosidase by site-directed mutagenesis and heterologous expression. *Journal of Biological Chemistry* **269**, 2283 -2291 (1994).
166. Nagy, J.K. & Sanders, C.R. Destabilizing mutations promote membrane protein misfolding. *Biochemistry* **43**, 19-25 (2004).
167. Yu, Z., Sawkar, A.R. & Kelly, J.W. Minireview: Pharmacologic chaperoning as a strategy to treat Gaucher disease. *FEBS Journal* **274**, 4944-4950 (2007).
168. Sawkar, A.R. et al. Chemical chaperones and permissive temperatures alter the cellular localization of gaucher disease associated glucocerebrosidase variants. *ACS Chemical Biology* **1**, 235-251 (2006).
169. Grabowski, G.A. & Hopkin, R.J. Enzyme therapy for lysosomal storage disease: Principles, practice, and prospects. *Annu. Rev. Genom. Human Genet.* **4**, 403-436 (2011).
170. Barton, N.W. et al. Replacement therapy for inherited enzyme deficiency - Macrophage-targeted glucocerebrosidase for Gaucher's disease. *New England Journal of Medicine* **324**, 1464-1470 (1991).
171. Platt, F.M. et al. Inhibition of substrate synthesis as a strategy for glycolipid lysosomal storage disease therapy. *Journal of Inherited Metabolic Disease* **24**, 275-290 (2001).

172. Lachmann, R.H. Miglustat Oxford GlycoSciences/Actelion. *Current Opinion in Investigational Drugs* **4**, 472-479 (2003).
173. Futerman, A.H., Sussman, J.L., Horowitz, M., Silman, I. & Zimran, A. New directions in the treatment of Gaucher disease. *Trends in Pharmacological Sciences* **25**, 147-151 (2004).
174. Cohen, F.E. & Kelly, J.W. Therapeutic approaches to protein-misfolding diseases. *Nature* **426**, 905-909 (2003).
175. Sawkar, A.R. et al. Chemical chaperones increase the cellular activity of N370S  $\beta$ -glucosidase: A therapeutic strategy for Gaucher disease. *Proceedings of the National Academy of Sciences of the United States of America* **99**, 15428 -15433 (2002).
176. Lin, H. et al. N-octyl-beta-valienamine up-regulates activity of F213I mutant beta-glucosidase in cultured cells: a potential chemical chaperone therapy for Gaucher disease. *Biochim. Biophys. Acta* **1689**, 219-228 (2004).
177. Steet, R.A. et al. The iminosugar isofagomine increases the activity of N370S mutant acid  $\beta$ -glucosidase in Gaucher fibroblasts by several mechanisms. *Proceedings of the National Academy of Sciences* **103**, 13813 -13818 (2006).
178. Kornhaber, G.J. et al. Isofagomine induced stabilization of glucocerebrosidase. *Chembiochem* **9**, 2643-2649 (2008).
179. Maegawa, G.H.B. et al. Identification and characterization of ambroxol as an enzyme enhancement agent for gaucher disease. *Journal of Biological Chemistry* **284**, 23502 -23516 (2009).
180. Desnick, R.J. & Schuchman, E.H. Enzyme replacement and enhancement therapies: lessons from lysosomal disorders. *Nat Rev Genet* **3**, 954-966 (2002).
181. Jain, A.N. Surflex: Fully automatic flexible molecular docking using a molecular similarity-based search engine. *Journal of Medicinal Chemistry* **46**, 499-511 (2003).
182. Almlöf, M., Brandsdal, B.O. & Åqvist, J. Binding affinity prediction with different force fields: Examination of the linear interaction energy method. *J. Comput. Chem.* **25**, 1242-1254 (2004).
183. Marelius, J., Kolmodin, K. & Åqvist, J. *Q Manual for Version 5*. (Uppsala University: 2004).
184. Jorgensen, W. & Rives, T. The OPLS potential functions for proteins - energy minimizations for crystals of cyclic-peptides and crambin. *J. Am. Chem. Soc.* **110**, 1657-1666 (1988).
185. *MacroModel*. (Schrodinger, LLC: 120 West 45th Street, 32nd Floor, New York, NY 10036-4041, USA, ).



186. King, G. & Warshel, A. A surface constrained all-atom solvent model for effective simulations of polar solutions. *J. Chem. Phys.* **91**, 3647 (1989).
187. Lee, F.S., Chu, Z.-T., Bolger, M.B. & Warshel, A. Calculations of antibody-antigen interactions: microscopic and semi-microscopic evaluation of the free energies of binding of phosphorylcholine analogs to McPC603. *Protein Engineering* **5**, 215 -228 (1992).
188. GROMACS 4.5 Online Reference. at <<http://manual.gromacs.org/>>